



Measuring affective states from technical debt: A psychoempirical software engineering experiment

Downloaded from: <https://research.chalmers.se>, 2023-05-05 03:22 UTC

Citation for the original published paper (version of record):

Olsson, J., Risfelt, E., Besker, T. et al (2021). Measuring affective states from technical debt: A psychoempirical software engineering experiment. *Empirical Software Engineering*, 26(5). <http://dx.doi.org/10.1007/s10664-021-09998-w>

N.B. When citing this work, cite the original published paper.



Measuring affective states from technical debt

A psychoempirical software engineering experiment

Jesper Olsson¹ · Erik Risfelt¹ · Terese Besker¹ · Antonio Martini² · Richard Torkar^{1,3}

Accepted: 7 June 2021 / Published online: 22 July 2021
© The Author(s) 2021

Abstract

Context Software engineering is a human activity. Despite this, human aspects are under-represented in technical debt research, perhaps because they are challenging to evaluate.

Objective This study's objective was to investigate the relationship between technical debt and affective states (feelings, emotions, and moods) from software practitioners.

Method Forty participants ($N = 40$) from twelve companies took part in a mixed-methods approach, consisting of a repeated-measures ($r = 5$) experiment ($n = 200$), a survey, and semi-structured interviews. From the qualitative data, it is clear that technical debt activates a substantial portion of the emotional spectrum and is psychologically taxing. Further, the practitioners' reactions to technical debt appear to fall in different levels of maturity.

Results The statistical analysis shows that different design smells (strong indicators of technical debt) negatively or positively impact affective states.

Conclusions We argue that human aspects in technical debt are important factors to consider, as they may result in, e.g., procrastination, apprehension, and burnout.

Keywords Technical Debt · Affective States · Software Development · Psychoempirical Software Engineering · Empirical Study · Bayesian statistical analysis

1 Introduction

Software engineering is very much a human activity, but this is sometimes forgotten. When proposing hypotheses, analyzing results, and discussing implications for the industry, we

Communicated by: Emerson Murphy-Hill

✉ Jesper Olsson
research@jesperolsson.se

¹ Department of Computer Science and Engineering, Chalmers and University of Gothenburg, SE-412 96, Göteborg, Sweden

² Department of Informatics, University of Oslo, N-0373, Oslo, Norway

³ Stellenbosch Institute for Advanced Study (STIAS), Wallenberg Research Centre at Stellenbosch University, Stellenbosch, South Africa

researchers sometimes neglect to factor in human aspects (Lenberg et al. 2015). So, too, is the case for technical debt research (except for a handful of studies on morale, e.g., (Besker et al. 2020)). This paper intends to amend this deficit and provide evidence showing that technical debt has noticeable adverse effects on software practitioners' feelings.

Technical Debt (TD) is a financial metaphor (Cunningham 1992), typically used within software engineering to explain long-term costs of short-term benefits (Ampatzoglou et al. 2015). It is a communicative aid for bridging the knowledge gap between software practitioners and business decision makers. Hence, if the metaphor was to miscount (or not account for) pivotal cost-benefit factors, the effect could be detrimental to software companies.

The current definition of TD was agreed upon during the 16162 Dagstuhl seminar (Avgeriou et al. 2016): "In software-intensive systems, technical debt is a collection of design or implementation constructs that are expedient in the short term, but set up a technical context that can make future changes more costly or impossible. Technical debt presents an actual or contingent liability whose impact is limited to internal system qualities, primarily maintainability and evolvability."

The definition is nuanced, incorporates decades of research, and offers a shared understanding of TD. Among many other things, it emphasizes that TD is a software development artifact in its own right and that TD acquisition is not necessarily intentional nor visible. A list of various consequences was also synthesized, but it fell short in recognizing the effects of TD on the human aspects of software engineering.

This paper aims to fill that gap by assessing five different design smells (proxies for design TD) to understand if, how, and why these smells impact participants' affective states during their development work.

In this study, we address this gap by employing a mixed-methods approach (including an experiment) and following guidelines for psychoempirical software engineering research ("research in software engineering with proper theory and measurement from psychology" (Graziotin et al. 2015c)). The study collected empirical data ($n = 200$ data points from $N = 40$ participants) on how design TD influences the so-called affective state of software practitioners. Applying Bayesian multi-level models revealed, among other findings, strong evidence that certain design smells (notably cyclic-dependencies) caused the subjects displeasure. The qualitative analysis suggests that many practitioners experience anxiety from high amounts of TD, and their responses vary along a maturity scale.

In more concrete terms, the research objective of this study is to investigate the relationship between TD and affective state from the point of view of software practitioners. This objective is supported by three research questions, which are listed below and further elaborated on in Section 3.

RQ1: How do software practitioners' affective state change in the presence of design smells?

RQ2: How do changes in affective state align with professional characteristics (e.g., formal education, work experience, or work context)?

RQ3: How do software practitioners reason about the relationship between affective states and technical debt?

The results of this study provide important insights and show that psychological factors also need to be acknowledged as a consequence of TD. The results show, for instance, that different kinds of design smells impact participants' affective states differently. When assessing how the affective state aligns with the practitioners' professional characteristics, the results show that work experience correlates with submissiveness. Lastly, practitioners reason, e.g., that negative affects often coincide with TD, but can be viewed as opportunities for improving the code base.

The sections of this paper are laid out as follows. The following section presents related work in the research areas of TD and human aspects of software engineering, individually and jointly. Section 3 describes the research design and methods employed. Next, Sections 4–5 present the quantitative and qualitative analyses, respectively. The study is discussed in Section 6, limitations and threats to validity are presented in Section 7, and the paper is concluded in Section 8.

2 Related Work

Much of the current literature on Technical Debt (TD) pays particular attention to technical or financial perspectives. This study breaks with such traditions to observe TD through the lens of human aspects of software engineering. Hence, for full appreciation, the reader should be familiar with the background of the two research fields.

Recounted firstly is previous research on TD in general. Appropriate nomenclature and central findings are outlined before introducing the specific type of TD investigated in this study. Secondly, we describe software engineering research on human behavior, emphasizing recent studies on the topic of feelings, emotions, and moods, and the recommendations concerning measurement instruments from psychology. One of those instruments, the Self-Assessment Manikin (SAM), was employed in this study and is explained in detail.

Once these two branches (i.e., the research area to be broadened, and the facet used to do so) have been covered, related work is listed. That is, existing research items that have used similar lenses and investigated challenges encountered in the TD literature. Those items are briefly reviewed to clarify how this study fits into the current body of knowledge.

2.1 Previous Research on Technical Debt

Technical Debt (TD) was conceptualized a few decades ago by Cunningham (1992) as a financial metaphor for how early misunderstandings of a problem domain might hamper future development unless the software is refactored to incorporate knowledge gained. Since then, the term has received much attention in both academia and industry. Today, the metaphor is widely used as a communicative aid for explaining internal software quality problems to non-technical stakeholders by emphasizing the extent to which the software must compromise its ability to meet the needs of the future to meet the needs of the present (Cunningham 1992; Avgeriou et al. 2016; Ampatzoglou et al. 2015; Fernández-Sánchez et al. 2017; Ernst et al. 2015).

One of the main strengths of TD is that much of its terminology originates from finance. As noted by Ampatzoglou et al. (2015), the two most commonly used terms in TD research are *principal* and *interest*, i.e., the cornerstones of financial debt. In software engineering, the former expresses the effort required to turn the current quality of some development artifact into its optimal level—the latter concerns how this sub-optimal level of quality leads to extra effort in later development iterations.

Several studies have shown that TD has significant negative consequences that can be detrimental to software companies (Tom et al. 2013; Li et al. 2015; Besker et al. 2018a; Ampatzoglou et al. 2015; Fernández-Sánchez et al. 2017). The interest does away with a substantial portion of development time (Besker et al. 2017; 2019), and may grow non-linearly if left unattended (Martini and Bosch 2017). Further, TD tracking and TD management are uncommon in the software industry, and when encountered, the processes are typically immature (Guo et al. 2011; Ernst et al. 2015; Martini et al. 2018a).

Despite its severity, TD is difficult or impossible to measure directly, and assessments typically rely on measurement proxies known as software smells, i.e., indicators of (internal) software quality issues (Fontana et al. 2017; Ganesh et al. 2013; Garcia et al. 2009; Sharma and Spinellis 2018). Naturally, empirical studies, such as this one, face the same issue when they need to exemplify TD items.

So far, we have outlined the previous research on TD in general, by giving an account of its history, terminology, and critical findings. The next paragraphs will focus on a type of TD known as Design TD (DTD), which our investigation is based on.

True to its name, DTD is TD found in software design, i.e., sub-optimal constructs in the software system's structure and behavior. As such, its boundary to, e.g., architectural TD (ATD), is disputed. Some researchers merge the two (Tom et al. 2013). Others separate them (Li et al. 2015; Alves et al. 2016) according to definitions that typically are too vague or subjective to form disjunct sets (Alves et al. 2014; Alves et al. 2016).

Such disagreements propagate to the categorization of software smells (Garcia et al. 2009), which results in some smells, e.g., cyclic dependencies and hub-like dependencies being considered either design smells (Ganesh et al. 2013) or architectural smells (Fontana et al. 2017).

To reduce the risk of misinterpretation, this study will not merge the two categories. The investigation is concerned with small, local problems, in isolated parts of the software system that can be comprehended easily. The findings should not be confused with the large concerns highlighted in recent ATD research, e.g., Ernst et al. (2015) and Besker et al. (2018a).

2.2 Previous Research on Human Aspects of Software Engineering

A growing body of literature recognizes the importance of interdisciplinary research between software engineering and psychology (Cruz et al. 2015). Both academia and the industry acknowledge that software engineering tasks are human activities and, thus, impacted by human aspects (Boehm and Papaccio 1988; Feldt et al. 2010; Colomo et al. 2010; Tamburri et al. 2013; Fagerholm et al. 2015).

For many years, such studies were dispersed, but in 2015 Behavioral Software Engineering (BSE) was proposed as a common platform for research concerned with “the study of cognitive, behavioral, and social aspects of software engineering performed by individuals, groups, or organizations” (Lenberg et al. 2015).

Out of the many tracks in this research area, one concerns *affective states* (or *affects*, for short), i.e., feelings, emotions, and moods. Previous studies have linked affects to, e.g., debugging performance (Khan et al. 2011), analytical ability (Graziotin et al. 2014), and productivity (Graziotin et al. 2015b).

This study is placed firmly within this track and is part of a sub-field called psychoempirical software engineering (PSE), i.e., software engineering studies that use theory and measurements from psychology (Graziotin et al. 2015c). This article follows the Graziotin et al. (2015c) guidelines for conducting PSE research.

According to these guidelines, this study's objective is best met by subscribing to the *dimensional framework* and employing the *Self-Assessment Manikin* (SAM) instrument for measuring affective states (Graziotin et al. 2015c). Within the dimensional framework, affects are expressed through several distinctive dimensions, e.g., the models represent affective states along three continua: pleasure–displeasure (valence), arousal–nonarousal (arousal), and dominance–submissiveness (dominance) (Graziotin et al. 2015c; Russell and Mehrabian 1977).

In more concrete terms, according to Graziotin et al. (2015b), these dimensions can be understood as follows. Valence is the attractiveness (or adverseness) of an event, object, or situation, while arousal is the intensity of emotional activation or the sensation of being mentally awake and reactive to stimuli. Finally, dominance is the sensation of control of the situation; one's skills are perceived to be higher than the challenge level for the task.

The recommended instrument, the SAM, measures affects through pictorial representations (Fig. 1) of the three dimensions of the models (Graziotin et al. 2015c; Lang 1980; Bradley and Lang 1994; Morris et al. 2002). Developed by Lang (1980), the instrument has, over the decades, been subjected to extensive validation research (Morris 1995) and seen used in numerous studies, see (Morris 1995; Betella and Verschure 2016).

According to Bradley and Lang (1994), the graphic design of the SAM has many benefits. The lack of verbal components means that the SAM can be administered to a broader population range, including individuals with a non-English mother tongue or language disorders, and children. Additionally, the SAM can measure direct affective reactions, as it can be filled out in a short amount of time and eliminates cognitive processing (Morris et al. 2002). Further, Morris (1995) argues that the use of stylized characters, as opposed to photographs of humans, makes the SAM less susceptible to many types of biases.

However, because SAM relies on self-reporting, the scores are not standardized according to objective reference points. Although individuals are consistent with themselves (within measurement), the ratings cannot be assumed to be consistent between individuals (between measurement) (Graziotin et al. 2015c). In other words, two individuals could rate

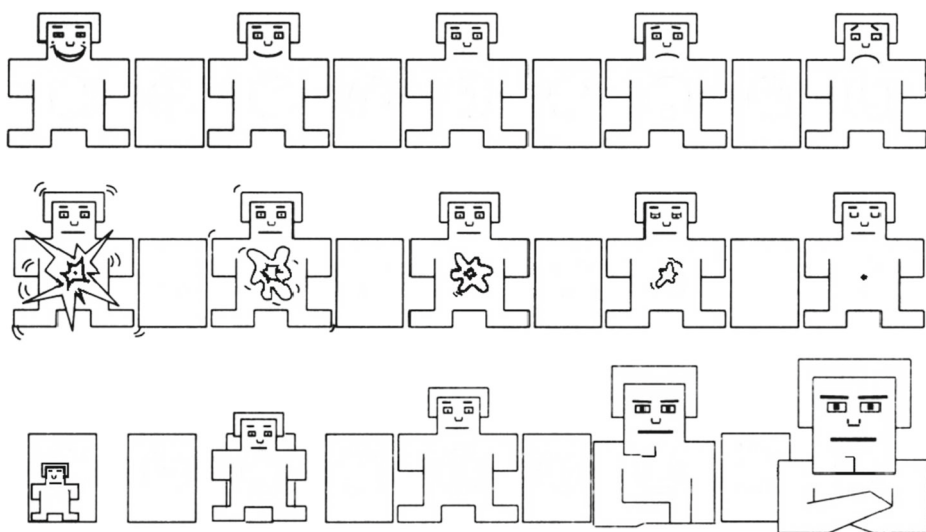


Fig. 1 The SAM measurement instrument. SELF ASSESSMENT MANIKIN ©Peter J. Lang 1994

the same affective state in two different ways. Consequently, investigations administering the SAM should follow a within-subject (or repeated measures) design (Graziotin et al. 2015c), which also follows the latest recommendations in general (see Gelman 2018).

Additionally, it is important to recognize that the SAM is not suited for all types of affective state research; Graziotin et al. (2015c) emphasize that the instrument is designed to measure affective reactions in response to a stimulus (in our case, design smells). For example, the SAM would be unfit for studies aiming to investigate how happy software practitioners are, generally (Graziotin et al. 2015c).

The SAM is protected by copyright law, but the instrument and instructions for proper administration (Lang et al. 1997) are available for non-profit academic research.¹

2.3 Interdisciplinary Research on TD and Human Aspects

Data from several secondary studies reveal that few TD studies have investigated the relationship between TD and human aspects (Tom et al. 2013; Li et al. 2015; Ampatzoglou et al. 2015; Alves et al. 2016; Fernández-Sánchez et al. 2017; Besker et al. 2018a). Rather, the predominant concerns have been technical and financial aspects, e.g., software quality or cost of future changes.

When human aspects are addressed in TD research, the most frequently investigated topic is morale. A negative correlation was proposed early by Tom et al. (2013) based on anecdotal evidence found in web blogs. Since then, empirical investigations have corroborated the connection, including previous articles of our own, see (Besker et al. 2020).

Spínola et al. (2013) performed a survey on TD folklore and found medium to high consensus among software practitioners that TD is related to their morale. In conjunction with interviews, a survey was also carried out by Besker et al. (2020) to determine how occurrence and management of TD affect developers' morale. Their findings show that the existence of TD negatively impacts morale, but also that morale is increased by proper TD management.

Although a common misconception, morale is not the same thing as affective state (Graziotin et al. 2015a; Peterson et al. 2008). Hence, to the best of our knowledge, there are no previous TD studies investigating affects and even fewer that directly measure how software practitioners respond to TD items.

In addition to morale, some empirical studies have offered evidence for TD harming the software practitioner's psychology. Lim et al. (2012) found that developers are more reluctant to incur TD because its consequences become a part of their daily work. Similarly, such reluctance may arise due to developers predicting that the sub-optimal construct needs to be corrected sooner or later, and that task would fall on them (Yli-Huumo et al. 2014). However, these findings were somewhat opportunistic and limited, as neither study set out with the research objective of investigating such questions.

TD research has thus far shown lukewarm interest in the relationship between TD and human aspects. However, the topic has also been approached from the PSE direction, and those studies present interesting empirical findings. Graziotin et al. (2017) surveyed software practitioners concerning causes for unhappiness, and established that low code quality and coding practices, and being stuck in problem-solving, were among the most significant factors. Additionally, in a later study, Graziotin et al. (2018) investigated the adverse effects

¹Information about how to obtain the SAM can be found at <https://cse.phhp.ufl.edu/Media.html>

of developer displeasure and found, among many other types of consequences, *lower code quality* and *discharging code* (extreme cases of productivity and quality drop, in the form of deleting parts of the code base).

Not only are these factors intimately connected with TD, but they pose the threat of vicious cycles: Low code quality causes unhappy developers, and unhappy developers produce low-quality code. Unfortunately, the studies did not drill down into this problem, which could answer questions such as its probability and severity. Nor was the issue approached specifically from the TD perspective. Clearly, our study differs from the previous PSE studies, as it seeks to investigate affects regarding specific TD items.

In conclusion, prior research shows that investigating human aspects concerning TD is a promising prospect. To manage TD more effectively, we need to understand how software practitioners, as human beings, can be factored into the trade-offs between short-term and long-term benefits. However, the current body of knowledge is limited, and both academia and the software industry would likely benefit from further clarification.

3 Methodology

As suggested in the previous section, our research topic has received little attention despite interesting initial findings. Consequently, the study design must acknowledge the limitations posed by such research gaps, e.g., validation against previous findings may be impossible.

One of the countermeasures implemented in our design is choosing a mixed-methods approach, i.e., collecting both quantitative and qualitative data. This decision is appropriate because it enables the study to improve validity, e.g., the results from one analysis could corroborate or rebut findings from the other. In this study, data were gathered from three sources: A repeated-measures experiment (quantitative), a questionnaire (quantitative), and a semi-structured interview (qualitative).

Another central countermeasure is the high transparency achieved by providing a replication package for this publication.² It contains complementary information and all material needed for reproducing the study, as it is infeasible to present all details within the scope of this article.

To demonstrate this study's overall study design, we have constructed a holistic research design model as illustrated in Fig. 2. As shown, this study was conducted in three different phases: a design, an execution, and a synthesis phase. The figure also illustrates the different performed activities within each phase and references the sections describing these activities. If more information exists in the replication package, this is also pointed out (using the tag `repl.pkg`).

3.1 Goals

This study seeks to examine the relationship between design smells and software practitioners' affective states. Thus, it tries to understand the importance of human aspects as a factor in TD. Among other things, we hope that the answers to our research questions will spark further interest in considering software practitioners when making trade-offs between

²<http://doi.org/10.5281/zenodo.4537801>

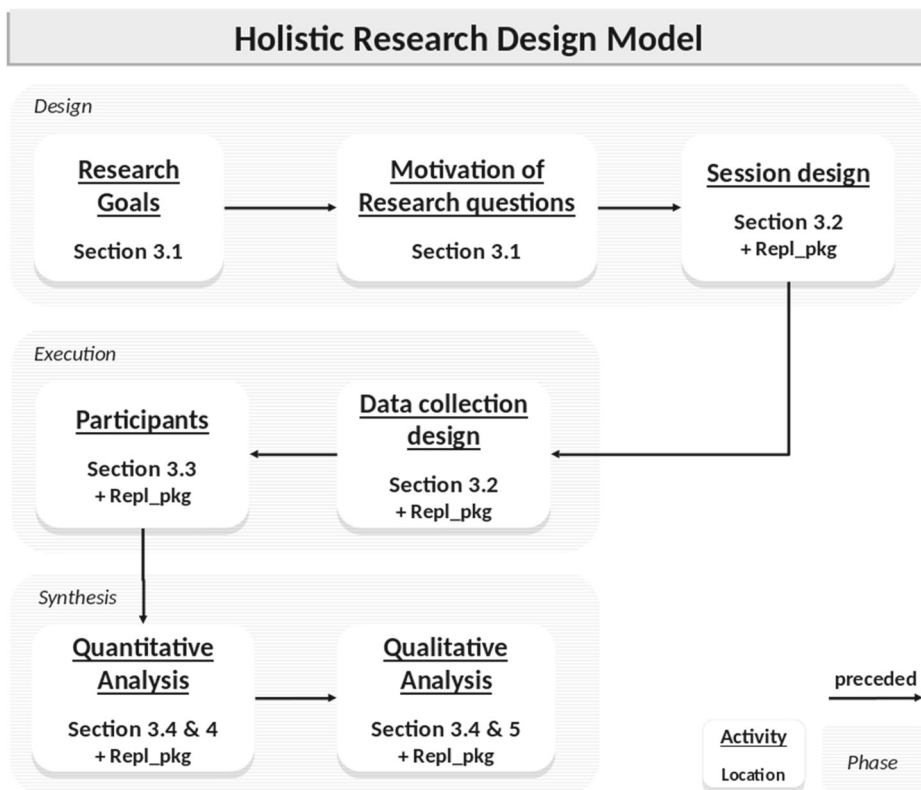


Fig. 2 The design of the study

short-term and long-term benefits. This goal begs for persuasive evidence, which can be provided through empirical research.

RQ1 (see Section 1) will be answered by conducting a within-subjects experiment. The data are analyzed via (Bayesian) statistical analysis: We employ dynamic Hamiltonian Monte Carlo to sample multi-level models. This research question aims to investigate the actual relationship between affects and DTD, without being colored by the participants' (nor the researchers') preconceived notions. As for delimitations, this RQ will examine a handful of design smells and consider affects from the presented models' perspective alone.

The motivation behind RQ2 is to see what role individual differences play. Because the study examines affects, the experimental units must be human participants, which opens up many exciting characteristics that could be studied. However, while data for various factors could be collected with ease, there are trade-offs to consider, e.g., transparency and confidentiality. Since the data are open (see the replication package), many characteristics that could easily identify an individual (e.g., gender or ethnicity) were not recorded.

Finally, RQ3 was included to understand the topic's appearance in the software industry. Hence, this research question is broader than the other two and of a more exploratory nature. Giving voice to the practitioners' reflections on affects and TD can increase understanding in a broader context and reveal peripheral issues.

3.2 Session Design

As this study collected three sets of data, its design is a substantial part of this article. Since there are many constructs to keep track of and clarify, we will use a few different viewpoints. The first viewpoint is that of *sessions* and is modeled in Fig. 3.

From this perspective, the study was designed as 90-minute sessions, one for each participant. At the start of their session, the participant received instructions (pre-task instructions) outlining the study and the session. The participant obtained these in three steps:

- 1) reading, understanding, and signing a document describing the treatment of, and their rights regarding, collected data (confidentiality assurance);
- 2) listening to instructions for, and seeing examples of, how to use the measurement instrument—which relies on self-reporting (SAM instructions); and
- 3) hearing a description of what activities they will perform during the experiment (task description).

Next, during the second part of the session (measurement sitting), quantitative data were collected from a repeated-measures experiment. For this part, as well, the participant went through three steps (please note that being of a repeated-measures design, the second and third steps were conducted five times):

- 1) using the measurement instrument on a practice task (anchor point);
- 2) pausing briefly (deacclimatization period); and
- 3) using the measurement instrument on a task (scenario).

In the last part (post-task interview), the two remaining data sets were gathered: quantitative data from a questionnaire and qualitative data from a semi-structured interview. These were presented to the participant in one step each:

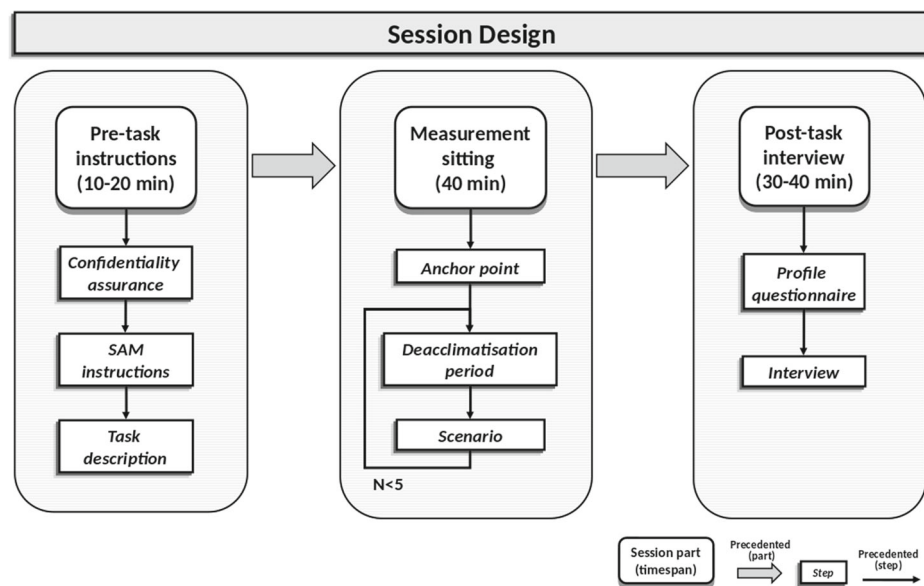


Fig. 3 The session view of the study: 90-minutes sessions, conceptually comprising three parts with eight steps

- 1) filling out answers to questions about their professional experience with software (profile questionnaire) and;
- 2) talking and answering questions about how they perceived the study and their view of code maintainability and feelings (interview).

Thus, the session perspective is concluded. This description has given an overview of what the participants did, between being greeted by the researchers to saying goodbye. It also introduced concepts that are key to understanding the study design but did so on a high abstraction level. Further details on these concepts can be found in the replication package.

Next, we consider the study from the perspective of *data collection*. Three sources of empirical data (experiment, questionnaire, and interview) were gathered from the participants. As shown in Fig. 4, each of these data sets was designed around one of the RQs, i.e., the experiment for RQ1, the questionnaire for RQ2, and the interview for RQ3. Similarly, the experiment data and the questionnaire were modeled in the same statistical analysis, while the interview data underwent thematic analysis.

First, the experiment set out to understand the relationship between affects and DTD. From this goal, it followed that, ideally, all factors except for the amount of design debt (explanatory variable), should remain constant. Then, what was measured was the participants' affective state in terms of valence, arousal, and dominance (response variables).

However, since the experiment was of the repeated-measures variety, its design was more complicated. While the explanatory variable still represented the amount of design debt, there was not one but five such variables (one for each repetition or *scenario*). In other words, as the participant progressed through the experiment, they would encounter five different scenarios: ScA, ScB, ScC, ScD, and ScE. Within each scenario, the participant received one treatment and then reported their affective state.

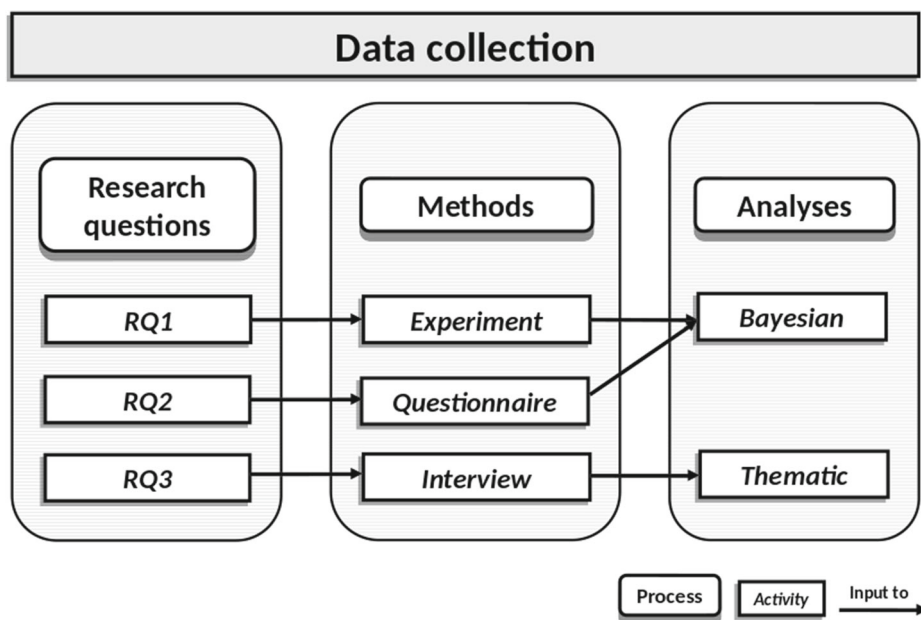


Fig. 4 The relationships between the RQs, methods, and analyses

Because design debt is difficult to measure, each response variable had two levels and represented whether its design smell (see Table 1) was present or had been refactored away. That is, the scenario variant where the smell had been removed had a lower (*L*) amount of technical debt than its partner variant (*H*).

The scenarios were derived from Suryanarayana et al. (2014), which in turn is based on Ganesh et al. (2013). Because smells are not necessarily indicative of definite quality problems (Sharma and Spinellis 2018), smell catalogs such as Garcia et al. (2009) were considered inappropriate for the experiment.

Moving on to the second method, the questionnaire aimed to investigate how professional characteristics factor into the participants' responses. The questions are listed in Table 2.

The third method, the interview, was designed to answer RQ3 and explore the topic of TD and human aspects beyond the delimitation of this study. Because the quantity of previous studies is limited, the study gains extra benefits from validating and contextualizing its findings. Hence, caution should be exercised when limiting the participants' divergent thinking and, thus, the data's richness. Therefore, the participants were not constrained to talk merely about DTD.

Instead, the participants were allowed to speak more or less freely about their perception of affects and software maintainability. The questions listed in Table 3 were asked at opportune times during the interview to light-handedly steer it. These were complemented by probing questions, i.e., follow-up questions to the participants' reasoning.

Because the interviews had a broader scope than this study, the thematic analysis used to answer RQ3 considered a subset (highlighted in green) of the interview questions, namely IQ4.1, IQ4.2, and IQ5.

Thus, the data perspective is concluded. It presented how the research questions can be traced to the selected methods and analyses. Further, the general structure of the methods was explained, including the questions asked of the participants.

The third and final perspective is the *materials* perspective, which is illustrated in Fig. 5. Their description is deferred to the replication package, where the experimental protocol also is included.

3.3 Sample

Forty software practitioners from 12 companies participated in this study. The participants were obtained through convenience sampling, but covered a diverse set of professional characteristics, e.g., their experience came from many different business domains (such as automotive, finance, and renewable energy) and ranged from 1 to 35 years. All participation was voluntary and based on informed consent and anonymity.

Table 1 The scenarios used in the experiment and the smells they embody

ID	Smell	Smell category
ScA	Missing Encapsulation	Encapsulation smell
ScB	Missing Hierarchy	Hierarchy smell
ScC	Broken Modularization	Modularization smell
ScD	Cyclically-Dependent Modularization	Modularization smell
ScE	Rebellious Hierarchy	Hierarchy smell

Table 2 The questionnaire

ID	Type	Description
Q1	Closed	My highest level of completed academic education is _____
Q2	Closed	My education major (e.g., computer science, electrical engineering, software engineering, ...) was _____
Q3	Closed	I have working experience with software for _____ years.
Q4	Closed	My current role (e.g., architect, developer, tester, ...) is _____
Q5	Closed	The programming language I am most experienced in is _____
Q6	Closed	My currently preferred programming language is _____
Q7	Closed	Most of my working experience comes from the following domain (e.g., telecom, healthcare, finance, ...) _____
Q8	Open	Do you have any additional comments concerning this questionnaire?

3.4 Analysis Procedure

Two different analyses were performed in this mixed-methods study. For the quantitative part, a Bayesian statistical model was implemented and executed in *R* (R Core Team 2020). The procedure is available in the replication package.³

The qualitative data was analyzed by following the guidelines for thematic analysis by Braun and Clarke (2006). Thematic analysis is frequently applied in both psychology (Braun and Clarke 2006) and software engineering (Cruzes and Dybå 2011).

The flexibility of thematic analyses stems from several choices that the researchers must make when deciding how to conduct the analysis (for a discussion about each choice's advantages and disadvantages, see (Braun and Clarke 2006)). For this study, the analysis was *inductive*, searched for *semantic themes* and theorized *essentialistically*. In other words, we coded the interview transcripts in a data-driven fashion without trying to fit them into a pre-existing coding frame. Themes were then identified and interpreted based on what was explicitly articulated within the data set.

The primary reason for these decisions is the small amount of previous research on the relationship between TD and the human aspects of software engineering. For example, the *inductive* approach does not rely on existing theory to the same extent as the *theoretical*. Similarly, it seemed more prudent to identify the themes at the *semantic* level, given the exploratory nature of this investigation. Otherwise, the likelihood of projecting personal beliefs onto *latent* themes could be excessive. The same reasoning underpinned the choice of performing an *essentialist* analysis. In particular, previous research on human aspects of TD did not seem to lend sufficient support for theorizing socio-cultural contexts and structural conditions (beyond little more than pure speculation), as is sought with the *constructionist* perspective.

Since the qualitative analysis aimed to discover the most central ideas and themes (rather than most, or all of them), the analysis's size was determined by salience rather than (thematic) saturation (Weller et al. 2018). This decision is somewhat uncommon in software engineering research, so a short motivation is in order.

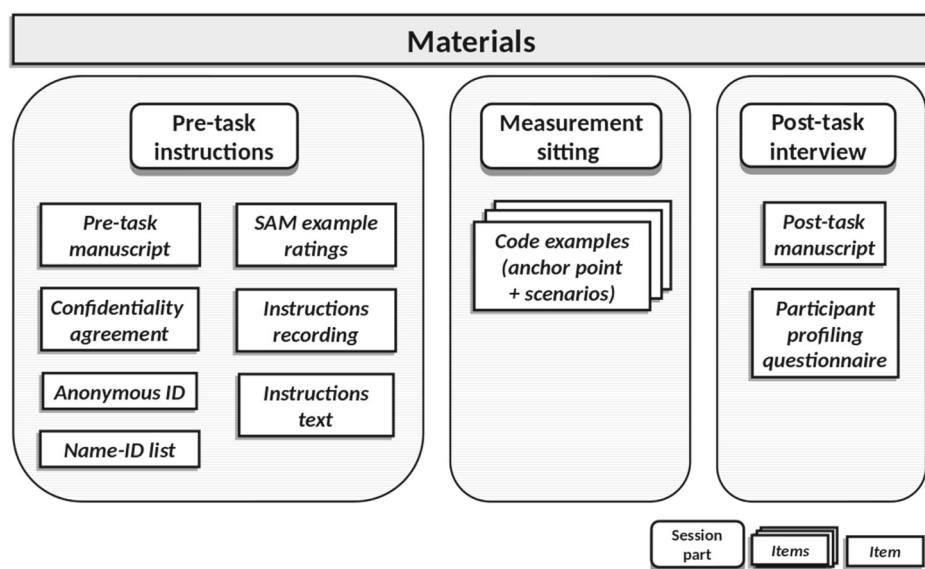
Salience is the idea of analyzing qualitative data regarding the most prominent items, and can be contrasted with saturation, i.e., until the set of all unique items is *believed* to have

³<http://doi.org/10.5281/zenodo.4537801>

Table 3 The common questions of the semi-structured interview. The thematic analysis used to answer RQ3 is considered a subset (highlighted in green) of the interview questions

ID	Type	Description
IQ1.1	Open	Could you please tell us more about your daily work. What type of tasks do you normally encounter?
IQ1.2	Open	How do those tasks make you feel?
IQ1.3	Closed	Do you face challenges in those tasks?
IQ1.4	Open	How do those challenges make you feel?
IQ1.5	Closed	Are those feelings frequent?
IQ2	Open	In contrast to challenging tasks, what sorts of feelings would you say you get from routine tasks?
IQ3	Closed	Do you think that anything outside of this experiment did impact your responses today?
IQ4.1	Open	Would you please tell us how you experienced the code examples?
IQ4.2	Open	What about the software design in the examples?
IQ5	Open	What would you say are the differences between the scenarios we provided and software one encounters in industry?
IQ6	Closed	Did you find SAM difficult to use or understand?
IQ7	Open	That was all of the questions that we had for you. Is there anything you would like to add?

been exhausted. For a broad range of research objectives, saturation would be superfluous, as salient items are, unsurprisingly, more prevalent and more culturally significant than non-salient items (Weller et al. 2018). In other words, many research questions can be answered with smaller sample sizes than what would be required to claim saturation.

**Fig. 5** The experimental materials used in the different parts of the session

The point at which thematic saturation is reached depends not only on the domain size, but also on the number of responses per person (Weller et al. 2018). Consequently, salience may be the more appropriate alternative when it is difficult to know the size of the domain or the set of ideas (Weller et al. 2018) (as is the case in this study).

At the same time, the importance of probing questions should not be overlooked: When the investigation aims to obtain most of the most important ideas and themes in a domain (as is frequently the case in qualitative research and particularly in open-ended interviews), a smaller sample with extensive probing is commonly more productive than a large sample with casual or no probing (Weller et al. 2018). Thus, salience should be used with caution, unless the data collection is designed with this in mind.

Because 10 interviews are sufficient to reliably capture up to 95 % of the most salient ideas (Weller et al. 2018), that number of data items was randomly selected for the data set (out of the 39 items in the interview data corpus).⁴ Indeed, this study's necessary sample size might be even lower, as we used probing techniques during the interviews, e.g., repeating phrases the interviewee uttered when working with the scenarios and asking for more information.

4 Quantitative Analysis and Results

Forty subjects participated in the experiment, and each subject contributed with five measurements to estimate our outcomes. Also, the following data were collected: Educational level (e.g., bachelor), the example used (the ten experimental artifacts, i.e., five artifacts in *L* and *H* setting), academic major (e.g., computer science), role (e.g., designer), language experience (e.g., Java), entities (i.e., level of complexity of the artifact), and years of work experience. The latter was scaled in order to improve sampling (i.e., $(x_i - \bar{x})/x_\sigma$).

Given the three outcomes valence, arousal, and dominance $\{V, A, D\}$, and the predictors listed above, the data consists of a matrix with 200 observations (rows) and 11 variables (columns), with no missing data.⁵

In this analysis, we employed Bayesian ordinal regression, using a cumulative model (for an introduction to Bayesian analysis, see Furia et al. (2019)). One could imagine two other potential models, i.e., the sequential model or the adjacent category model. However, since Likert (1–9) scales were used for the outcome, cumulative models are more suitable, i.e., the sequential model would be suitable if we want to analyze the number of correct designs predicted from experience. In contrast, the adjacent category model would be appropriate if we want to predict the number of correctly solved sub-items of a complex task—none of this was of interest to us (Bürkner and Vuorre 2019).

Several models were designed, and their relative out-of-sample prediction capabilities were evaluated iteratively. The final model, below, includes all relevant predictors and has the same out-of-sample capabilities as other comparable models. For model comparison, we used state-of-the-art model evaluation (Vehtari et al. 2017).⁶

Next, follows the design of the final model and the corresponding priors. If we want to make a comparison with a frequentist approach, then one could claim that we have fixed and

⁴A single participant asked not to be recorded during the interview and could thus not be included.

⁵The dataset, with analysis scripts and a Docker image, can be found at <http://doi.org/10.5281/zenodo.4537801>. R 4.0.2, rstan 2.21.2, and brms 2.13.9 was used for the analysis (R Core Team 2020; Bürkner 2017; 2018; Stan Development Team 2020)

⁶Pareto $k < 0.5$ and LOOIC = 2406.0.

random effects in our model (i.e., a mixed-effects model); however, in a Bayesian setting, we use the term multilevel model, since that allows us also to employ hyperparameters with corresponding priors.

$$V_i, A_i, D_i \sim \text{Cumulative}(\phi_i, \kappa) \quad (1)$$

$$\phi_i \sim \beta_1 \text{EDUCATION}_i + \beta_2 \text{EXAMPLE}_i + \beta_3 \text{MAJOR}_i + \beta_4 \text{ROLE}_i \quad (2)$$

$$+ \beta_5 \text{LANGUAGE}_i + \beta_6 \text{ENTITIES}_i + \beta_7 \text{EXPERIENCE}_i \quad (3)$$

$$+ \beta_{\text{SUBJECT}[i]} \quad (4)$$

$$\beta_1 \sim \text{Dirichlet}(2, 2, 2, 2, 2) \quad (5)$$

$$\beta_{\text{SUBJECT}} \sim \text{Half-Cauchy}(0, 2) \quad (6)$$

$$\beta_2, \dots, \beta_7 \sim \text{Normal}(0, 0.5) \quad (7)$$

$$\kappa \sim \text{Normal}(0, 5) \quad (8)$$

In the first line we model each outcome, $\{V, A, D\}$, using a cumulative likelihood. The parameters ϕ and κ are the linear regression and the intercepts, respectively, which we model for each outcome (i.e., we have eight intercepts for each outcome since the outcome was Likert scale 1–9).

In the next three lines, we have the linear regression. We have eight parameters we want to estimate, one for each of our predictors. The parameters β_1 and $\beta_{\text{SUBJECT}[i]}$ are special as we will see next.

On Line 5, we assign β_1 a Dirichlet prior. The Dirichlet prior is the multivariate generalization of the Beta distribution (a distribution commonly used to model a probability $[0, 1]$). Using Dirichlet, we can model an array of probabilities; i.e., in this case, we model five probabilities and use a very weak prior (the 2s), indicating that we do not have any prior knowledge. The reason we use a Dirichlet here is monotonicity, i.e., the predictor EDUCATION is an ordered categorical variable indicating the level of education. We, thus, want to model the probability separately for each of the categories in education.

Continuing on Line 6 we assign β_{SUBJECT} a Half-Cauchy(0, 2) prior. This prior is common when modeling standard deviations and allows only positive real numbers (\mathbb{R}^+). To analyze variability in this way goes by many names, e.g., random effects or varying intercepts. The reason we use it is due to our following the latest recommendations by designing the experiment to collect within-person measurements (Leek et al. 2017), i.e., each subject has been randomly allocated several tasks and, thus, we model the variability of each subject to partially pool the estimates, to avoid overfitting.

Proceeding to Line 7, we assign the priors Normal(0, 0.5) for the remaining parameters while, on the last line, we assign the prior Normal(0, 5) to all intercepts for each outcome. (It is common to assign a broader prior for intercepts.)

The careful reader would react to what could be perceived as tight priors for several parameters, i.e., Normal(0, 0.5). However, first, using Normal(0, 0.5) on six parameters still makes an impressive standard deviation, $(6 * 0.5)^2 = 9$, and, second, the combination of all priors established a *nearly uniform prior* on the probability scale, i.e., prior predictive checks and a sensitivity analysis were conducted.

Since we used dynamic Hamiltonian Monte Carlo to sample, we also have several diagnostics. In our case, the model showed no indications of a biased posterior, and diagnostics (\hat{R} , effective sample size, and trace plots) indicated that the chains had converged. Posterior predictive checks showed that the data swamped the priors (see Fig. 6a and b for a visualization of the prior predictive checks and posterior predictive checks).

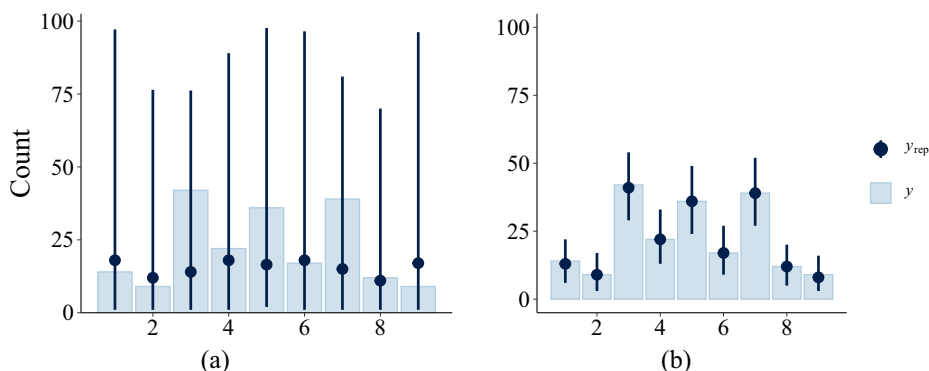


Fig. 6 Prior and posterior predictive checks (y is the empirical data, and y_{rep} are 100 draws from the prior (a) and posterior (b) probability distributions). The left plot shows the prior predictive checks (where no empirical data was used). The uncertainty is considerable (the lines), and the median values (the dots) are approximately the same for all items on the Likert scale, like it should be since only the priors are used. Compare this to the right plot, where we have drawn samples from the posterior probability distribution, i.e., we have fitted our model with data, the data has provided evidence, and thus the priors have been what is commonly referred to as ‘swamped’, since the uncertainty has decreased

Continuing this section, we will next look at the output from the model. First, we will present the standard deviations for each outcome’s random effects and any interesting population-level effects. Then, we will predict outcomes while fixating specific parameters. The final part will present the results of the hypothesis testing (Bayes factor).

Analyzing the variance, there is not much difference in the uncertainty of the estimates concerning σ for our three outcomes, as the standard deviations’ credible interval mass vary from 0.88 (σ_V) to 1.1 (σ_A). In short, the uncertainty for each outcome, $\{V, A, D\}$, is very much the same, but, notably, valence (V), has the lowest standard deviation $\sigma = 0.39$, while arousal (A) has the largest standard deviation, $\sigma = 0.87$, indicating more uncertainty in between-subjects variability. This can be interpreted as that the within-subject design and analysis we employed was beneficial (it was important to model different dispersions).

Analyzing the estimates, and their corresponding 95% credible intervals, led to 5 estimates being singled out as interesting (Table 4). Four were significant on the arbitrary 95%-level (i.e., not crossing zero), while one is strongly positive, albeit not significant on the 95%-level.

Since Experience has much probability mass on one side of zero ($[-0.05; 0.56]$), we will analyze it further to understand its predictive ability better. Before we analyze Experience further, let us look at the role Entities (i.e., the complexity of each task) has on the outcome. If it is not positive, then one could argue that they have had the wrong effect.

Table 4 Parameters of interest

Outcome	Parameter	Est.	Est. Error	l-95% CI	u-95% CI
Dominance (D)	EXAMPLE (BL)	−0.78	0.34	−1.43	−0.12
Valence (V)	EXAMPLE (BH)	0.73	0.34	0.07	1.39
Valence (V)	EXAMPLE (DL)	−0.83	0.35	−1.52	−0.14
Valence (V)	EXAMPLE (CL)	0.72	0.36	0.02	1.42
Valence (V)	EXPERIENCE	0.25	0.16	−0.05	0.56

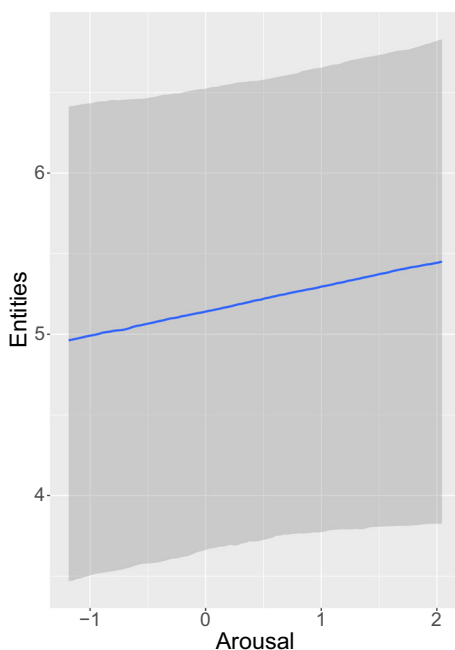
To investigate Entities we need to determine what covariate values to use. One possible way to do this is to set all values to their mean for continuous variables, while the reference category is used for factors, and then examine the conditional probabilities our posterior probability distribution provides us with. In Fig. 7, we see a positive trend, which indicates that the model has been able to capture the role that complexity plays correctly.

Finally, we would like to see the role Experience plays by analyzing it more carefully. If we turn our attention to Fig. 8a–c, we see that the role it plays differs, depending on our outcome. For Valence (V), we have a positive effect, i.e., the more experienced the subject, the higher the response on the Likert scale, while the opposite holds for Arousal (A) and Dominance (D). Here, it is crucial to keep in mind the direction of the SAM, i.e., an increase in V score means more displeasure; arousal increases as A decreases; low D scores denote submissiveness.

Having analyzed the conditional effects, we now turn our attention to measuring the strength of the evidence we have gathered. Our tests will *not* examine the significant population-level effects, which we list in Table 4; after all, we know that they are significant on the traditional 95%-level. Instead, we will focus on the contrasts between Low (L) and High (H) settings for our predictor Example. This means that we can present the results as several hypothesis tests (5 artifacts times 3 outcomes equals 15 tests in total). Since we have a posterior probability distribution, we do not have to, generally speaking, worry about multiple tests, which is often the case in a frequentist setting (Gelman and Tuerlinckx 2000; Gelman et al. 2012).

For hypothesis testing, we will use Bayes factor to avoid the usage of p -values and, thus, to receive verdicts both in favor of and against a given hypothesis (Goodman 1999a; 1999b). For our accept/ reject decisions, we follow recommended practices as presented in Table 5 (Kruschke 2010).

Fig. 7 Conditional effect of Entities in the model. The more complex an entity (i.e., the more to the right we move on the x -axis), the higher the outcome on the Likert scale (y -axis). In this case, we looked at the outcome A (arousal), but the same trend is visible in all three outcomes. The x -axis has been scaled, with 0 corresponding to median complexity. (The line is the median outcome, while the gray area is the 95% uncertainty around the median)



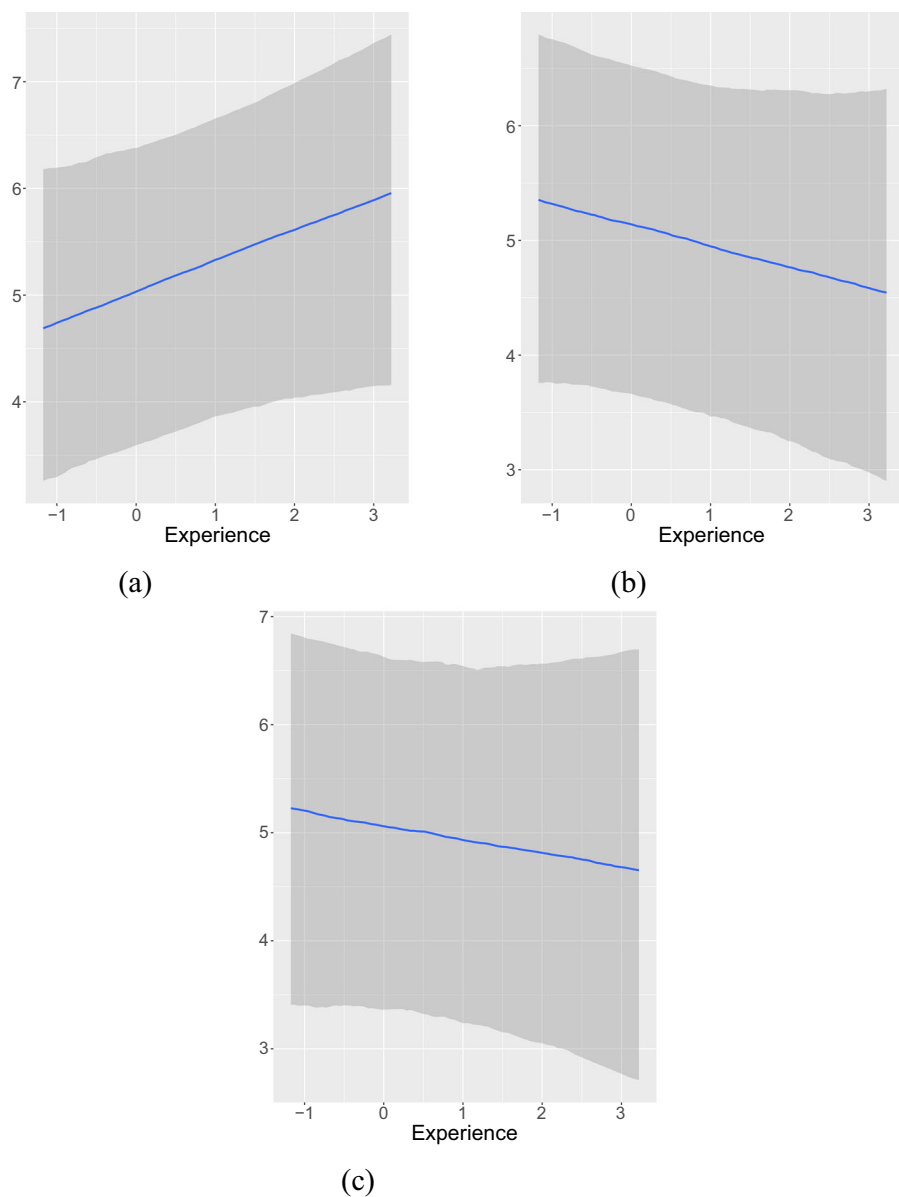


Fig. 8 An overview of the conditional effects on Experience, given our three outcomes $\{V, A, D\}$. Lines correspond to the median, while the gray area is the 95% credible interval. For valence (V), we have a positive effect, i.e., the more experienced the subject, the higher the response on the Likert scale, while the opposite holds for Arousal (A) and Dominance (D)

Our hypothesis tests were unidirectional and, thus, tested that Low < High, e.g.,

$$H_0 : \text{Example}_{\text{AL}} < \text{Example}_{\text{AH}},$$

which is to be interpreted as Example A Low is less than Example A High (and we analyze this inequality for each of our outcomes $\{V, A, D\}$).

If we plot the posterior probability distributions for each hypothesis test (15 in total), one can perhaps better see what a ‘significant’ effect means in the context (Fig. 9a and c).

4.1 Effect Sizes

Looking at Fig. 9a–c one sees three hypotheses that indicate strong evidence, i.e., Examples *D*, *A*, *C* in outcome *V* (valence). In the two former cases, we have evidence for H_1 , while in the latter case we have evidence for H_0 . Analyzing the effect sizes for these results is wanted. However, we also see two more results that could potentially also be of interest.

In Fig. 9c, one can see that there are some probability distributions classified as providing *moderate* evidence for H_1 or H_0 , respectively (but they are still fairly close to a quantile). These are Examples *B*, *C*, and *D*. Even though we do not have strong evidence speaking in favor (or not) of a hypothesis, it could be of interest to see what this entails concerning effect size.

In short, we want to see, on average, how large an effect size it would be to move from *H* to *L* for each of the six Examples. By drawing samples from our posterior probability distribution, we can easily compare the difference between levels. We leave all variables according to what we have in the sample (e.g., the distribution concerning Experience is the same) and vary only the Example level to see what this means on the outcome scale. Table 6 provides us with an overview of the six effect sizes.

One can conclude this section by claiming that we have some interesting effects, some even based on substantial evidence. These are summarized in the box below as findings F1–F11.

Findings for RQ1:

- F1 Cyclically-dependent modularization (ScD-H) is less pleasant than its refactored (ScD-L) counterpart (strong evidence).
- F2 Missing encapsulation (ScA-H) is less pleasant than its refactored (ScA-L) counterpart (strong evidence).
- F3 Broken modularization (ScC-H) is more pleasant than its refactored (ScC-L) counterpart (strong evidence).
- F4 Missing Hierarchy (ScB-H) is, likely, less dominating than its refactored (ScB-L) counterpart (moderate evidence).
- F5 Broken modularization (ScC-H) is, likely, less dominating than its refactored (ScC-L) counterpart (moderate evidence).
- F6 Cyclically-dependent modularization (ScD-H) is, likely, more dominating than its refactored (ScD-L) counterpart (moderate evidence).

Findings for RQ2:

- F7 Work experience, likely, correlates with submissiveness (moderate evidence).

Additional findings:

- F8 Refactored Missing Hierarchy (ScB-L) yielded particularly submissive responses.
- F9 Missing Hierarchy (ScB-H) yielded particularly displeasing responses.
- F10 Refactored Cyclically-Dependent Modularization (ScD-L) yielded particularly pleasing responses.
- F11 Refactored Broken Modularization (ScC-L) yielded particularly displeasing responses.

Table 5 Decision thresholds for hypothesis testing using Bayes factor, according to Kruschke (2010)

Symbol	Evidence ratio	Description
**	> 10	Strong evidence for H_1
*	3–10	Moderate evidence for H_1
?	1–3	Anecdotal evidence for H_1
?	$1/3$ –1	Anecdotal evidence for H_0
*	$1/30$ – $1/10$	Moderate evidence for H_0
**	$< 1/10$	Strong evidence for H_0

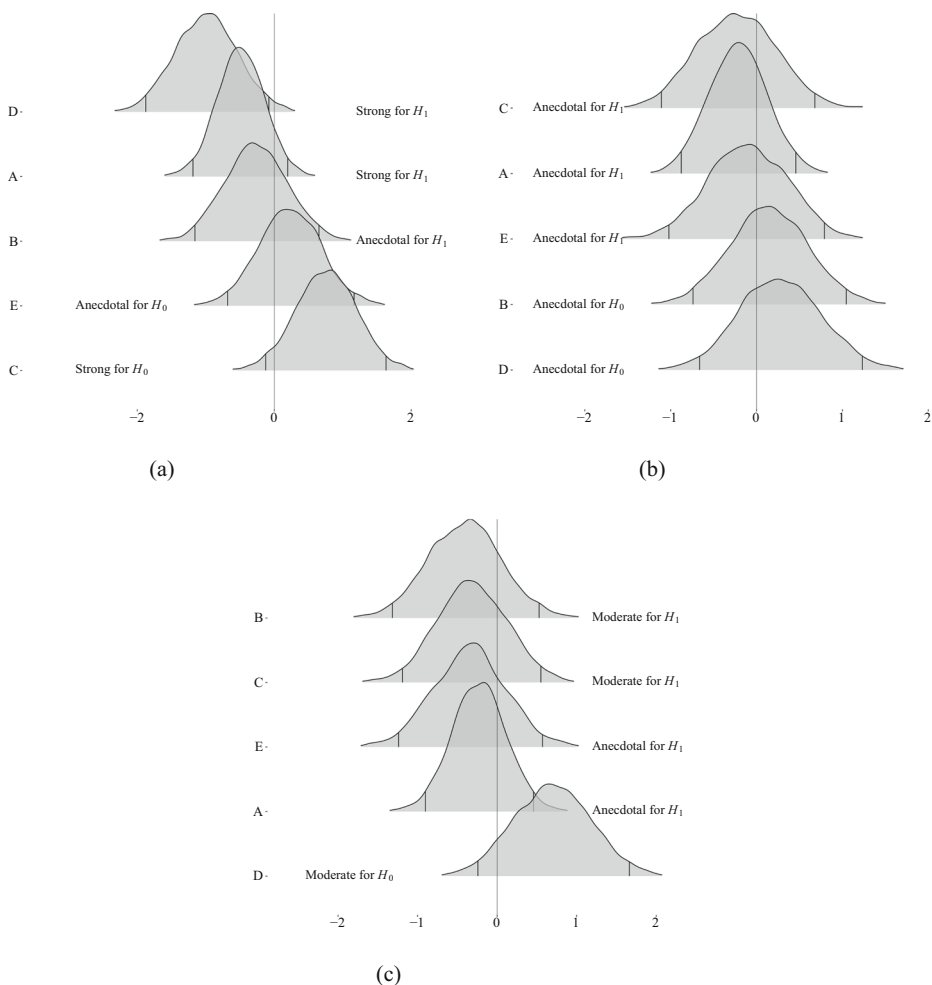


Fig. 9 A visual overview of all hypothesis tests, given our three outcomes $\{V, A, D\}$ (x -axis is the contrast). On the y -axis, Examples (A–D) are ordered according to the direction of evidence *starting with the most negative direction*. Next to each distribution, a short note clarifies the results of the tests (according to Table 5). Finally, the distributions have 2.5% and 97.5% quantiles drawn in the tails. As an example, artifacts D, A, and C, in outcome V (valence) indicate strong evidence. In the two former cases we have *strong* evidence for H_1 , while in the latter case we have strong evidence for H_0

Table 6 Raw effect sizes from posterior samples (10,000 draws) of the posterior predictive distribution. These samples have higher variance than samples of the means of the posterior predictive distribution since residual error is incorporated. The first three rows present raw effect sizes where the hypothesis test found strong evidence, while the last three rows show where there was moderate evidence. The median column is the size of the effect (on the outcome scale) for the contrasts $L-H$. If we look at the first row we see an effect size of -1.0 , i.e., the difference between Low–High, for Outcome V and Example D , is -1.0 on the Likert scale with the quantiles $[-1.5, 1.8]$. This should not be confused with the hypothesis tests we conducted (Fig. 9a and c), which tested if $\text{Low} < \text{High}$

Outcome	Example	Min.	1st quant.	Median	3rd quant.	Max.
Valence (V)	D	-3.6	-1.5	-1.0	1.8	3.9
Valence (V)	A	-3.0	-1.0	-0.5	0.0	2.0
Valence (V)	C	-1.6	0.2	0.8	1.3	3.8
Dominance (D)	B	-3.2	-1.0	-0.4	0.1	2.7
Dominance (D)	C	-3.1	-0.9	-0.3	0.2	2.6
Dominance (D)	D	-2.5	0.2	0.78	1.3	3.8

5 Qualitative Analysis and Results

Analyzing the data set (which predominantly concerned the participants' general experience of TD, rather than the experiment scenarios) revealed that the participants have strong and negative affects toward TD and are inclined to talk about their reactions. Their argumentation was clearly of the stimulus-response variety, i.e., they viewed TD as an action they are exposed to, leading to counteractions. The participants' discussions centered around what one might think of as defense or coping mechanisms for said stimulus.⁷

The thematic map (including two themes and five sub-themes) constructed during the analysis is included in Fig. 10. The first theme (three sub-themes) describes the participants' reflections (with regard to affective state) on undergoing TD intense areas (*Undergoing TD*), e.g., encountering TD, when working with some other task.

Among its sub-themes, we first consider *Procrastination*. At its core, this sub-theme is about instances where practitioners try to delay or avoid dealing with the debt or its consequences. Often, this is related to the sense of feeling overwhelmed when facing TD.

Procrastination may surface in several different forms. For example, one interviewee reported that TD could cause task abandonment. “the more, like, bad code I see [in the same place], the more, like, bored and [indifference] [...] It's like, ‘[vocabulary of quitting], I give up’. It's like, ‘it's too much now, I give up.’”

This feeling of resignation was echoed by another practitioner, who also suggests that tightly coupled code is cognitively taxing. “it had this instance of bit that implies that it knows about something else, so then you have to start knowing about two places at once, in parallel, and that usually gets super messy. [vocabulary of distaste] Yeah, so it's, sort of, being in control and being able to fix it.”

At the same time, *Procrastination* is not constrained to low levels of arousal. Quite the opposite, in some instances, it can lead to an impulsive and risky overhaul of parts of the codebase: “I would throw away and rewrite it”.

⁷These are established terms within psychology, and the surrounding theory could not be delved into for the scope of this study. In this article, we will instead use the term *psychological rebound* to avoid overloading the terms.

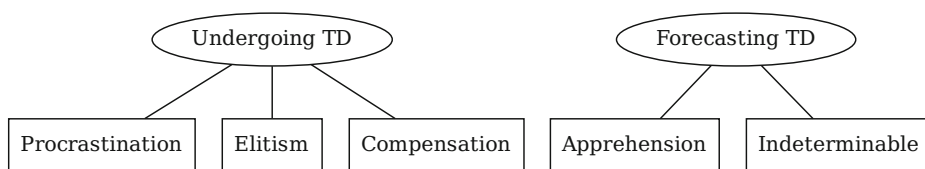


Fig. 10 Thematic map of how software practitioners reason about TD in tandem with affects

From these examples, it is clear that TD can cause *psychological rebound* effects that are harmful to the software project in ways that go far beyond the human aspects perspective. For example, abandoning tasks because of TD can upset backlog prioritization or result in project slippage. Similarly, the urge to overhaul the code base could, e.g., invalidate prior trade-off analyses.

Unsurprisingly, the participants were aware of the consequences and severity of *Procrastination*. One interviewee said, “*I think the detrimental part is when you feel like you don’t wanna touch it [...] even if I do touch it in the end, it will take a longer time before I actually dare.*”

Next, the second sub-theme is *Elitism*. It encompasses reactions to TD, violating some expectations that one holds oneself, one’s colleagues, or the code base to. In the case of *Elitism*, these expectations typically do not represent a shared set of values and beliefs among the parties. Hence, the discourse in this sub-theme was notably flavored by negative interpersonal dynamics.

Elitism is reflected in several different affects that appear to fall on a wide scale of blaming the author of the code. One example of a low amount of blame was one participant who expressed disappointment. “*if you have a great design, a great architecture, following the SOLID principles. That are loosely coupled. [Then,] they [code problems] are easy to fix. The problems. Easy to change. That is the most important, to me. So there are some—they are fundamentals of how I think when I design a program. So [code] violating those principles make me feel very sad.*”

As can be seen, this suggests that the code base itself influenced the participant’s affects, i.e., more or less decoupled from its author. On the other hand, another interviewee, who experienced distrust, accentuates the author’s (perceived) skill and does not separate it from the quality of the code: “*I’ve seen things where people mix really bad indentation, combined with not having, like, opening and closing brackets for **for**-loops, for example. Using, like, short notation. We can have, like, one-liners after **if**-statements, for example. I mean, those things are just terrible, ’cause you don’t know what belongs where. It’s messy and there are, like, no, like, blank lines between—additional spacing between things or anything. It’s just a bunch of code, with wrong indentation. Sometimes indented, sometimes not. And unclear what belongs to which statements. [...] it’s easier to spot it [than architecture]. And it’s so, like, something I really think people should know how to do. It’s so basic, in programming. So, yeah, I think so. It makes me a bit more worried, so to say, when I see that stuff. ’Cause it’s very much easier to do correctly.*”

Continuing on this blame scale, examples arise where the code is de-empathized in favor of focusing on its author. For instance, one interviewee expressed scorn and a notion of coding style reflecting one’s personality. “*I get a bit annoyed with people that try to be too smart with the programming language. They know, like, a short way of writing things, and they know exactly what happens. [...] So I’m more for, like, writing simple, easy to*

understand code. So that everyone that follows, can easily make changes to it. So yeah, that annoys me a bit: when people try to be too clever. They wanna show off that they are smart, by using, like, weird functions of a language.”

Viewed together, these examples suggest that *Elitism* may arise from the misalignment of quality expectations. However, this is perhaps not obvious to the practitioners, as the focus is not on addressing these alignment problems. Clearly, *Elitism* threatens to cause conflict among employees, but can also rationalize TD by acknowledging the debt as a result of business constraints: *“However, I also feel that when I read someone else’s code, that’s really bad—or shit, or something—I also realize that this might have been done under pressure, depending on the project and stuff. So I accept these technical debts better. Unless it’s just plain bad and not time-saving at all.”*

Elitism can be dangerous also when it does not result in (external) conflict. One interviewee highlights the risk of it causing high levels of stress. *“Yeah, this was people that were sort of in the more, like, architect roles, usually. Then they put on too much work on their shoulders. They were the guys that always wanted to do everything by themselves. And, sort of, tended to burn out after a while, ’cause they just had too much to do. You could see that they were stressed about it [soft deadlines].”*

The final sub-theme of the first theme is *Compensation*,⁸ which concerns constructively addressing TD. Often, the TD items are viewed as opportunities for improvement. As one interviewee put it, *“So, there definitely is this scope for improvement, but I would not call anybody else’s code as poor. [...] I generally do not get any negative feelings about it [code clones]. But I do look at it as an opportunity to improve the code myself.”*

Compensation is not limited to correcting an instance of TD, but it can encourage preventive actions, e.g., informing the code author about their mistake: *“Personally, I would use `git blame` to see who wrote the code and then, if I can contact them, I say ‘okay, next time, you should do it better. Because this, like, it may take a lot of time for others to trace their issues.’”*

One interviewee even suggested that affects can be leveraged to improve the code base, as they can act as software quality proxies: *“emotions aren’t bad or good. If a team member is that mad about something, I just use that as an indicator that something is bad in the code. So that person is right to be angry, and we can use that to either fix it, or use that as an argument for—like, in the future—like, let’s refactor this in the next sprint, or whatever.”*

Together, these extracts show that *Compensation* is related to TD management and, more specifically, tactics for addressing TD maturely or constructively. Please note how these tactics are concerned with the practitioners’ dominance concerning the code base. As one participant said, *“I want to rewrite it [code with inheritance issues]. [...] to improve it and to just, yeah, maybe so I don’t feel stressed about it. So I have control.”*

So far, we have presented the components of the first theme. Before continuing with the next theme, the interactions between these components should be analyzed. Note how all three sub-themes appear to be *psychological rebounds* for TD, albeit as different manifestations. *Procrastination* looks like an impulsive and naive, almost childish reaction to TD, where the practitioner does not acknowledge the consequences of their actions. These traits can also be seen in *Elitism*, but with regard to collaboration and teamwork rather than how the debt itself is approached. On the other hand, *Compensation* appears to be a manifestation of thoughtful consideration of how to manage the TD.

⁸In the behavioral (not the financial) sense of the word.

The second theme (two sub-themes) describes the participants' reasoning (with regard to affective state) when forecasting the consequences of TD (*Forecasting TD*), e.g., the effects of leaving TD unaddressed.

Its first sub-theme is *Apprehension*, which includes the anxiety of expecting future maintainability issues. A significant part of this sub-theme constitutes the participants' concerns about the extra psychological toll caused by TD. This toll can emerge when the practitioner believes there is a risk of the code leading to system failure. As one participant said, *"The spontaneous feeling was a bit stress about too much stuff going on. Too many components, and some strange dependencies. And too much inheritance. [...] Why [do I feel stressed]? 'Cause I can see myself maintaining that code. And I can see that code breaking in the long term. [...] 'Cause I don't want the system to break."*

This kind of uncertainty was echoed by another interviewee, who emphasized the toll of unforeseen consequences (ripple effects): *"for me, it comes back to, like, the control. I know that if I'm gonna touch this, I'm gonna pull a string, and then there's gonna come, like, a spider web with a spider in it. [...] You know that when you do something here, it's gonna affect something else."*

However, *Apprehension* is not limited to the technical considerations, as the psychological toll can also appear in the presence of tight schedules. As one practitioner put it, *"If you have time pressure to do something, and then you also know that you're in—I mean, 'this is gonna be hard to test. And to deliver it in time is gonna be tough.' Then it's super stressful. But if you don't have that pressure again, then it's easier again."*

Clearly, *Apprehension* is found in situations where the practitioner's dominance is on the submissive part of the scale, where they have low confidence in the code. Further, the extracts suggest that work tasks and business considerations are difficult to separate from their affective states. As one participant said, *"I mean, they [the technical and emotional viewpoints] are connected somehow. But through my years—my experience—I see a lot of problems with code violating these [SOLID design] principles. And that causes frustration when you try to fix bugs, improving the code, extend the code. So, it's more from a technical perspective, but they cause negative emotions."*

The last sub-theme is *Indeterminable*, which encompasses the difficulty of decoding TD. That is, understanding or sharing one's understanding of the TD in the system appears to be a non-trivial matter, which could play a key role in assigning value to TD items.

In industry, TD items are sometimes so opaque that professionals may not recognize them until they have paid a significant amount of interest. As one interviewee said, *"one time I was just gonna write some test for a thing we did. Then I realized the whole thing was such a debt-cluster that I just had to throw it away. I spent like three, four hours trying to help my team out. I didn't realize I did zero value [laughs] with that time."*

At the same time, TD might be widespread in the system, becoming a sort of background noise challenging to pinpoint. As one participant put it, *"sometimes you actually encounter some area that makes you really unhappy to be in. But then you also have these overarching stuff, that isn't really bothering you that much. But you always know it. You know it's always there. So it's way—it's less tangible. I would say it's, like, hard to identify. Hard to measure."*

Further, practitioners may recall areas with a high amount of TD but are sometimes unable or unwilling to articulate the problem constructively: *"I hear about '[vocal of complaint] this shitty part of the system.'"*

These extracts tell us that software practitioners have trouble estimating and communicating the consequences of existing TD items. However, as suggested by one interviewee, they may hold strong intuitions. *"In industry, it's more 'I know something is wrong. It feels like things are spread out like this. I just can't put my finger on it.' [...] The feeling I have*

in industry is more like, 'I know I'm gonna work in this area. I know it's gonna be horrible. I don't know what's gonna happen, exactly. Something is gonna show up. It's gonna take longer time. I can't give you a real estimate for how much it is to fix all of it, and I can't give you a real business value, because I just know it's gonna be hell.'

In conclusion, the analysis reveals that affects are very much a key aspect of TD. They provide an insight into the underlying mechanics for how software practitioners respond to TD items. These *psychological rebounds* may be a necessary consequence of TD and should not be ignored. The findings are further summarized in the following data extract and the box below (as findings F12–F24). *"if it's [the debt is] manageable or if I feel I can fix it, then it feels a bit okay. It's like, 'oh, this is a crappy thing someone did, but—whatever, it's fixable' in contrast to, like, 'this is just a nest of—we just need to re-engineer.' That makes you just angry inside. [...] you can definitely feel when it's '[vocabulary of excitement]', I can refactor this' or [...] [vocabulary of quitting], this is such a mess. I hate going into this code. I can't fix a bug here, 'cause there's just going to pop up things in other places. So, it's a mix. Depends on how much impact you can have on it, I think. Because it can be really fun to actually fix stuff. But when you can't, then it's like '[vocabulary of annoyance], angry.'*

Findings for RQ3:

- F12 Software practitioners experience (strong) affects from TD along all three dimensions.
- F13 When faced with high (overwhelming) levels of TD, practitioners will be reluctant to perform their work tasks.
- F14 Time pressure is sometimes a catalyst for negative affects.
- F15 Viewing TD items as opportunities for improvement appears to correlate with dominance toward the code base.
- F16 TD anxiety relates to code dependencies, ripple effects, and (the risk of) defect introduction.
- F17 TD anxiety appears to be correlated with submissiveness toward the code base.
- F18 Displeasure plays an important role in recognizing the presence and severity of TD.
- F19 Software practitioners sometimes get positive affects from amortizing TD.
- F20 Profanity frequently emerges in TD discussions.

Additional findings:

- F21 Quality processes sometimes get disrupted by software practitioners' affects.
- F22 Misalignment of quality expectations may result in interpersonal conflicts or burnout.
- F23 TD is challenging to decode (recognize, estimate, and communicate).
- F24 Violations of something the software practitioner considers fundamental appears to result in stronger affects.

6 Discussion

In this section, we tie together the results from the quantitative and qualitative analyses. We will continuously refer to the main findings (F1–F11 and F12–F24) at the end of Sections 4–5. First, we discuss the quantitative results related to the various scenarios and smells; to better explain the results, we then explore the quotes from the qualitative data occurring in correspondence with the analyzed smells. This allows us to explain how our results answer RQ1, or else how the smells influence the participants' affects. We compile a ranked list of which smells seem to have more impact.

We also discuss how changes in affective state align with professional characteristics (RQ2). We then take a broader scope and reason on the exploratory results from the qualitative analysis and what relationships we have found between affective states and technical debt (RQ3).

6.1 Case A: Missing Encapsulation

The quantitative analysis strongly suggests (F2) that the presence of the smell related to missing encapsulation in the code (ScA-H) causes the software practitioners to feel less pleasure (valence). This entails that practitioners consider the presence of this smell with disapproval rather than with indifference.

We do not seem to find other significant evidence related to the other two dimensions (arousal and dominance), which could imply that the practitioners do not consider this smell exceedingly threatening. This is also mentioned in the qualitative data, as one of the participants mentioned: *"Like, some of them were quite [vocabulary of annoyance] as solutions, but didn't really impact me that much. Like, the rectangle whatever—PNG-things [reference to ScA-H]. Like, yeah, I can refactor this in an afternoon."*

On the other hand, such a lack of strong feelings could be caused by the limited size and localization of the example and how easy it is to estimate the practitioner's refactoring. One of the interviewees mentioned *"So, it's—this, like, rectangle-PNG-thing [reference to ScA-H]—it's, like, I can really point to it. Show it. I can give an estimate for how much time is left and how much impact it is. The feeling I have in industry is more like, 'I know I'm gonna work in this area, I know it's gonna be horrible. I don't know what's gonna happen, exactly.'"*

In conclusion, the smell is recognized as a problem, but not as a high-priority one. Suppose we consider the strength of the resulting feelings and the participants' insights for this smell. In that case, we can conclude that the presence of this smell, although frowned upon, is perhaps not considered detrimental by the software practitioners.

6.2 Case B: Missing Hierarchy

We did not find evidence to support the hypothesis that this smell generates any negative feeling in the software practitioners. Surprisingly, on the contrary, we found (moderate) evidence (F4) that practitioners felt more dominant (dominance) in working with the code containing the smell (ScB-H).

On the other hand, this scenario was mentioned a lot in the qualitative analysis in rather negative terms. However, those comments often referred to the whole code and not to the specific smell. Although, at the same time, some participants explicitly mentioned the smell and suggested the correct refactoring. *"Yeah, one example, that had the private class there [referring to ScB-H], and that one I didn't like [...] Yeah, overall the checking of types in*

code like that: I think it's a sign of bad architecture, most of the time, when you have to check the type of objects coming in. Then you can probably—yeah, like I said—interface it out. And an interface coming in and you have the method on the interface and, yeah.”

The scenario itself could explain these seemingly contradictory results: debt-intense areas could bias the practitioner to distrust suitable constructs. After all, understanding and implementing a correct hierarchy involves a greater ability of abstraction. In other words, given that the original developers fell short in performing more straightforward tasks, confidence would be low to succeed in more demanding activities. As one of the interviewees said, in a different context, *“it's so, like, something I really think people should know how to do. It's so basic, in programming. So, yeah, I think so. It makes me a bit more worried, so to say, when I see that stuff. 'Cause it's very much easier to do correctly.”*

Another, less plausible, explanation would be that practitioners feel submissive (intimidated) *because* of the necessary abstraction skills, i.e., are not comfortable with such constructs. Here, it is essential to note the gap between recognizing a suspicious programming language construct (instanceof) and intimately understanding which abstraction would be suitable. The former is a low-level pattern detected by static code analysis (or even text search), while the latter often requires domain knowledge and experience. However, this seems unlikely, as most participants were experienced in object-oriented programming languages.

Finally, a third explanation could be that abstractions (by definition) remove details from the context. In other words, while beneficial for the system's maintainability, abstractions might, locally, result in less insight and, hence, less control.

In conclusion, the findings suggest that the smell is considered a problem despite its positive impact on dominance. The surrounding code's quality appears to confound individual TD items, but this effect needs to be verified in future studies.

6.3 Case C: Broken Modularization

Similar to Case B, we did not find evidence to support the hypothesis that this smell generates any negative feelings in the practitioners. However, we did find (strong evidence, F3) that they felt more pleasure (valence) and (moderate evidence) more in-control (dominance) when working with the code containing the smell (ScC-H).

This can be explained by the fact that the broken modularization smell consists of a widely recognized correct approach (modularization) applied in the wrong way (broken). In particular, the code that was modularized did not need to be (it consists of just variable declarations), and it should have been contained in the same abstraction (ScC-L). However, the participants' feelings might have been triggered by the presence of a better visual structure in ScC-H. The lack of a counter-effect for the displacement of the modularized code (ScC-H) could have different implications:

- 1) The positive feelings in the presence of modularized code far outperforms the negative feelings related to the sub-optimal use of such mechanism. This is also supported by one of the participants: *“It's, like, I could sort of see what had happened, I think. Like, the last one [reference to ScC-L] with the weird device. It looked like a container of data and someone plonked helper methods in it, maybe, I don't know. It's, like, I can see how that happened. I can move them without changing anything, so there won't be any ripple effects and I can still improve the code, for instance.”*
- 2) The practitioners could have overlooked the specific code that was modularized in an additional class, focusing more on the structure rather than on the code itself.

Alternatively, the participants might have thought that the additional class, the results of the modularization (which contains only variable declarations in our example), could have contained additional methods that were not displayed in our snippet.

- 3) While modularization is a well-known good practice, broken modularization is a less well-known bad practice among the practitioners. This can also be related to the language used by the participants and their familiarity with object-oriented programming. However, we did not find any evidence in the quantitative analysis supporting such an explanation.

In conclusion, we could not find evidence that this smell generates negative feelings in software practitioners. On the contrary, it seems as though the code with the smell was liked more, probably because the participants did not recognize (consciously or unconsciously) the misuse of modularization as significantly impacting.

6.4 Case D: Cyclic Dependencies

This is the smell for which we have quite strong evidence (F1) supporting hypotheses from literature (Martini et al. 2018b; Al-Mutawa et al. 2014). We can see how, for valence, the software practitioners reported extra-pleasure in the presence of code that is refactored (ScD-L), while at the same time, we register strong evidence that such code is much better liked than the one with the smell (ScD-H). Our analysis also reports moderate evidence for dominance (F6), where practitioners feel much more in control of refactored code (ScD-L) than the code containing the smell (ScD-H).

Despite such strong results, the smell was not often or explicitly mentioned in the qualitative answers of the practitioners dealing with ScD-H, if not for the two quotes below, which can be related to this specific smell: *“The spontaneous feeling was a bit stress about too much stuff going on. Too many components. And some strange dependencies.”* and *“the code doesn’t have to be perfect, or there could be some problems with the code. But if you have a great design, a great architecture, following the SOLID principles. That are loosely coupled. They are easy to fix.”* This could have happened because other smells or scenarios were more interesting to discuss, either because this example was not considered too challenging (perhaps because of the limited size of the example) or, possibly but perhaps less likely, because it was more noticeable and therefore less interesting.

In conclusion, we can consider this as evidence that the presence of cyclic dependencies generates stronger negative feelings in practitioners along at least two dimensions (valence and dominance).

Although this can be somewhat expected (cyclic dependencies is a well-known smell, probably more than the other smells), it is interesting to note how the degree of negative feelings for this smell far exceeds other smells. We find this plausible: Cyclic dependencies is the smell that tends to involve multiple entities (usually classes), which can generate ripple effects across the code. Also, the example that we propose here consists of just one dependency. In contrast, dependencies, especially if involving several entities, can become less noticeable and not so visible if they are not explicitly investigated, as shown in other publications, see Martini et al. (2018b) and Al-Mutawa et al. (2014).

6.5 Case E: Rebellious Hierarchy

We did not find even moderate evidence that this smell would generate either positive or negative feelings in the participants concerning any of the dimensions. Therefore, it is

difficult to draw conclusions on this smell and its impact on the practitioners' affects. In general, it seems as if the participants would be quite indifferent to one of the other proposed solutions. For example, even when encountering ScE-L, one of the practitioners mentioned: *"you had the Document [reference to ScE-L], yeah that was the wrong structure of the abstract class, I think, because you had all those methods, but only some of the implementation used. They didn't represent the same object. If you looked in the implementation, they had different actions or abilities. I think the public part of the implementations should be the same."*

The lack of evidence in itself combined with the quotation could point to three possible conclusions:

- 1) This smell is not considered a problem by practitioners, and it does not affect them.
- 2) Our example was not a good representation of the actual issue. Unfortunately, we did not find an existing implementation that would suit our experiment, so we had to adapt our snippet from Suryanarayana et al. (2014), removing any domain-specific reference that our participants would not understand. This process might have excessively simplified the smell.
- 3) The smell consisted of one short method out of ten, distributed in four (ScE-H) and five (ScE-L) classes, respectively. This could mean that the participants might have overlooked it in the time allowed for the task, perhaps focusing their attention on other code features, such as its structure (as mentioned in the previous quotation from a participant).

In conclusion, we cannot draw many conclusions from these results, although we can speculate that the participants have not recognized this as an issue affecting their feelings.

6.6 Comparison Across the Smells

We have so far reported our reflections, based on available evidence, on how and why the different smells have impacted the participants' affects. However, can we say something more about how the different smells compare to each other?

We report a summary (see Table 7), in which we compile a prioritized list of the smells based on the reported evidence. The smells are ordered by their negative impact on the software practitioners' affects. We also report if other impacts have been found on the refactored solution. To clarify the results, we have arranged the relationship positive/negative concerning the quantitative analysis, i.e., we do *not* consider the direction of the SAM. For example, in Table 7, 'negative impact on valence' means displeasure. We also highlight the strength and type of evidence supporting our conclusions.

6.7 Suggestions for Practitioners

Based on Table 7, we can undoubtedly suggest practitioners pay attention, especially to cyclic dependencies and missing encapsulations. As for the latter one, the practitioners mention that its refactoring would not be costly, which could make it a good candidate for a mandatory cleanup of the code before release.

Less vigorously, we also suggest that practitioners keep their eyes out for missing hierarchies. While the presence of this smell *increased* dominance that could be indicative of other problems, e.g., a need for domain knowledge acquisition, practitioners should consider taking action if developers start introducing this smell despite them recognizing it as bad practice.

Table 7 Prioritized smells according to our findings

Smell	Impact on affects	Dimensions	Other considerations
D—Cyclic dependencies	Negative Negative (likely)	Valence Dominance	Not explicitly discussed qualitatively
A—Missing encapsulation	Negative	Valence	Easy to estimate a refactoring
B—Missing hierarchy	Positive (likely)	Dominance	Recognized as bad practice, but overshadowed by the scenario code
E—Rebellious hierarchy	-	-	Seems to not have been recognized as an issue
C—Broken modularization	Positive Positive (likely)	Valence Dominance	Modularization seems to give positive feelings even if misused

As for the rebellious hierarchy, the study results do not allow us to draw firm conclusions, maybe because such smells might not be considered so upsetting by the participants. Finally, and probably surprising, it seems that *a modularization that is not entirely correct is not considered problematic*. At the same time, developers might tend to consider the pleasure of modularized code (even if containing a smell) better than smell-free code, which might be less modularized.

However, we need to notice that, for rebellious hierarchy and parts of broken modularization, the conclusions cannot be considered very strong, as our evidence is moderate. These results are also related to the influence of the smells on the participants' affects and do not consider other negative or positive effects. However, we consider our findings important to report, as developer unhappiness has been linked to harmful consequences (see Section 2.3). Further, the results should not be confused with the actual extra-maintenance effect these smells have in practice, although the two variables are most probably correlated.

6.8 The Effect of Experience

Our quantitative results do not point to correlations between the participants' professional characteristics and how the affective states changed. The most striking results are related to the experience of the respondents.

Experience has shown (moderate evidence, F7) to have a negative impact on Dominance. In other words, the more work experience the subject had, the more submissiveness they report.

A possible explanation for the increased submissiveness is that more experienced practitioners have dealt with the technical debt related to the smells for a longer time than junior ones. This may be caused by the fact that they have witnessed more of the technical debt's long-term negative impact, which may trigger additional caution for the smells.

These results also seem in line with our qualitative findings, primarily related to the maturity (see next section) with which practitioners undergo TD: we could argue that, with less experience (and, probably, less maturity), practitioners seem to want to ignore TD and avoid to worry about it (as highlighted by the *Procrastination* theme), hence the presence of lower submissiveness. Then, moving to a more elitist and compensating attitude toward TD as they gain more experience, they become additionally worried when encountering TD (as shown in the feeling of lacking control, mentioned in relation to the *Apprehension* theme).

Also, these results are in line with what is reported concerning how startup teams are composed and their inclination to incur TD. Besker et al. (2018b) report on interviewees from startups mentioning how it can be considered better to include a large part of junior developers in the initial team to make sure that TD is accrued (saving resources in the face of an initial high risk of failure). Experienced practitioners (apart from a small initial fraction) would be more suited for the growth and mature phase of a startup when TD needs to be removed before it becomes disruptive.

6.9 The Overall Effect of TD on Affects

Participants report that TD items activate a substantial portion of the emotional spectrum (three dimensions, F12), including vivid ones (e.g., profanity occurred, F20). Still, our experiment showed nothing concerning the arousal dimension. A plausible explanation is that participating in the experiment represents a different situation than encountering technical debt in real projects. The technical debt encountered during the experiment is not directly and negatively impacting the practitioners with, for example, extra-effort or additional bugs. This means that the arousal dimension could be triggered in a different context.

Many participants receive satisfaction from improving code (F15, F19). Mainly, being able to perceive their work as impactful causes pleasure. On the contrary, the uncertainty caused by code affected by TD and the consequent distrust in the code base are sources of negative feelings (F13). Architectural TD is considered a common source of negative feelings, especially for problems related to ripple effects (as, for example, in case *D* for cycling dependencies, F16).

Then the question is: why is TD so present in the software industry, and why is, e.g., code not continuously refactored?

First, as practitioners reported, stress is prevalent in the software industry. Several participants see deadlines as negatively affecting themselves and the product (F14). Avoiding TD requires more time, which would increase the stress in the presence of a deadline. This might mean that practitioners, to avoid stress, prefer to incur TD. Second, the participants mentioned that TD problems encountered in their daily lives are more extensive and more obscured than those in the experiment.

Another point of consideration was raised in the qualitative analysis, namely, that each sub-theme for undergoing TD is a *psychological rebound*. Further, there seems to be a sort of progression to them, which we will refer to as *maturity*, as we can draw parallels to our previous experience with group development models (Gren et al. 2017).

First, *Procrastination* (“Forming”) can be interpreted as a mechanism with little interest in improving the situation. Consequently, the practitioner will not attempt to share the team’s burden or attempt to shield its members from the harmful stimulus. Second, *Elitism* (“Storming”) involves questioning the code base and the *modus operandi*, which can be destructive and socially taxing unless adequately managed. Finally (note the absence of “Norming”) *Compensation* (“Performing”) illustrates a successful transition from defensive reactions to coping ones, with the participants focusing on facing up to the TD item and resolving it constructively.

6.10 Comparison to Related Work

In reviewing the literature, very little was found on the relationship between DTD and affective states, but several studies that investigated adjacent topics have found intriguing results. Many of the consequences of developer unhappiness demonstrated by Graziotin et al. (2018)

were echoed in our qualitative findings, e.g., *mental unease or disorder* (F13, F16, F17, F22). Perhaps most importantly, the quantitative analysis found a counterpart to *lower code quality*: Some design smells elicited an unhappier response among the participants than did other such smells (F1, F2). This could indicate a vicious cycle where TD leads to more unhappiness, which in turn leads to more TD.

Our previous paper on the relationship between TD and morale, Besker et al. (2020), found that TD negatively impacts morale, but also that morale is increased by proper TD management. These results are corroborated by F13, F16, F17, and F22; and F15 and F19, respectively. Hence, we provide further evidence in favor of the long-held belief that morale and TD are intertwined (Tom et al. 2013; Spínola et al. 2013).

The findings in this paper also corroborate several of those of Lim et al. (2012), e.g., developers fearing certain parts of the code base (F12, F13, F16, F17) and TD being difficult to communicate (F20, F22, F23, F24).

6.11 Implication for Research and Industry

In this study, we conducted an empirical investigation that joined the fields of TD and PSE. As demonstrated by Graziotin et al. (2018), affective states have important consequences for software engineering activities, and our findings provide solid evidence that design smells interlink with affective states. Accordingly, we present the argument that TD management should start factoring the human psyche into the decision-making processes.

Our findings, by themselves, constitute a compelling case, but do not stand alone. Although the human aspect is still a deficit area in the TD research, the combined results of Besker et al. (2020), Spínola et al. (2013), Lim et al. (2012), Yli-Huumo et al. (2014), and Tom et al. (2013)—many of which are corroborated by this paper—provide convincing grounds for our argument. Hence, we call for the research community to expand on the conceptual model of TD (Avgeriou et al. 2016). Figure 11 contains our proposition for how this part of the body of knowledge should be incorporated in our shared understanding of TD: The psychological factor would be explicitly acknowledged as a consequence of TD items.

The reason for expanding the model is to nuance how the research community and the industry view TD. The fact that software engineering is a human activity is often overlooked in investigations. While we recognize the challenge in putting a value on such aspects, that does not mean that we can turn a blind eye to their actual costs and benefits. Recognizing the psychological factors in TD will serve as the starting point for discussions and help the community converge on key concepts.

In particular, we encourage software engineers in the industry to engage in introspection, especially concerning stress and burnout. As surfaced in the qualitative data, many professionals face a psychologically taxing work environment, and until the consequences of their experiences are better understood, we advise caution. From our own experiences, the digital work environment, partly constituent of the code base, is seldom (if at all) regarded in analyses of occupational safety, health, and welfare.

Another crucial goal of TD management is to prioritize the removal of debt items that generate the worst current and future negative effects (or else, to use the metaphor, they have high interest attached). As repeatedly reported in the literature, this is a very difficult task, as measuring such interest is challenging and evidence is scarce. Measuring the affective states of practitioners in relation to different TD issues can be used as a proxy for such interest, or can at least provide additional insights on which items are perceived as the most “dangerous.” In Table 7, we provide a concrete example of how different smells (representing TD) impacted the participants’ affective states differently, which suggests a ranking across the

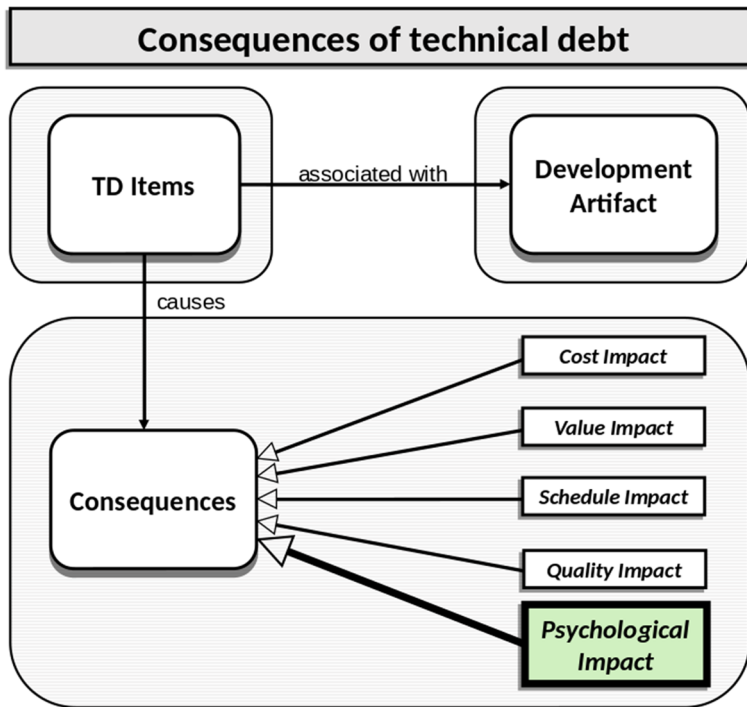


Fig. 11 Partial view of a conceptual model of technical debt, adopted from Avgeriou et al. (2016). The model has been extended to include the *psychological impact* in *consequences*

smells. In conclusion, a comprehensive catalog of smells and their impact on practitioners' affects could highly benefit the software engineering community.

Finally, we would like to address the psychological perspective in the context of education. Software engineering is difficult work and perhaps especially so because of the flexibility of the medium. Unlike other constructs, a code base is largely unconstrained by natural forces and can thus be perpetuated to unfathomable complexity. Unsurprisingly, resignation, frustration and even hate ensue. Many universities may want to consider explicitly educating their students about this reality and train them in how to engineer under such conditions.⁹ We argue that not imparting this knowledge would be an important oversight and urge the institutes to reflect on how our society is affected by the practitioner's emotional intelligence.

7 Limitations and Threats to Validity

Conducting empirical studies with human subjects is often a complex issue (Miller 2008), as it is often the case that the context is noisy and investigated effects are small (Gelman 2018). As such, the potential threats to validity are often numerous, and it would be infeasible to

⁹Project courses might *not* be the optimal choice as the work situation likely differs from that in the industry, but more research is needed to determine this.

discuss them all fully. This section presents what we consider the most significant validity threats to this study and the measures taken to mitigate them. The threats are categorized according to the aspects suggested by Wohlin et al. (2012) and Runeson and Höst (2009).

7.1 Construct Validity

This study set out with the aim of determining how DTD relates to the affects of software practitioners. One of the methods used was a repeated-measures experiment, where the participants were presented with five software design smells and their respective refactored versions. However, those smells were instantiated in code examples that originated from the same source, namely (Suryanarayana et al. 2014), which is not a scientific publication (although a derivative of one). Also, the source's purpose differs from ours in that the smells and the refactored versions are intended to be contrasted with each other. As previously stated, the rationale of this choice was a perceived deficit of suitable DTD representations in the research literature (for our purposes, at least).

These characteristics introduce threats to validity, but because they were identified before the data collection, countermeasures could be introduced. Since we were unable to find a way to eliminate these threats, we chose to monitor them and investigate the issue *post hoc*. This was done by introducing validity-checking questions (see Table 3, questions IQ4.1, IQ4.2, and IQ5) to the interview questions and analyzing the answers. (Further, IQ3, IQ6, and IQ7 checked other types of validity.)

In response to IQ5, the participants reported that the scenarios were representative of industry code, albeit atypically small and isolated examples of TD encountered in practice. Additionally, they confirmed that industry code would have impacted their affects to a greater extent. This suggests that the treatment was suitable, but also that the resulting data is an *underestimation* of the software industry's situation.

7.2 Internal Validity

The laboratory experiment part of this study was a repeated-measures design. While this approach lowers the threat of confounders because each subject's peculiarities are accounted for, it is more susceptible to learning effects.

Several countermeasures¹⁰ were taken to reduce learning effects. First, the participants were acquainted with the situation during the first phase of the session (pre-task instructions), and each received the same instructions for how to use the measurement instrument and their task. Second, the participants obtained practical experience with the procedure before the measurements (anchor point). Third, intermissions were used between measurements (deacclimatization periods) to lower the probability of any affects induced by previous scenarios carried over to later ones. Fourth, each participant was randomly allocated to one of two complementary treatment patterns designed to minimize bias. Fifth, and perhaps most importantly, *the order of the scenarios* was randomized for each participant.

7.3 External Validity

For this study, it is worth noting that the field of psychology is experiencing a replication crisis. Unfortunately, we have been unable to find consensus on concrete best practices for ensuring replicability and have instead chosen to adopt some propositions.

¹⁰Detailed in the replication package.

We have made our work as transparent as possible (see the replication package), under the constraints set by confidentiality, anonymity, and copyright. This includes the statistical analysis, the data, the procedure, and the experiment material.

Another issue concerned with external validity is the sampling strategy. In this study, we employed convenience sampling (further detailed in the replication package). The approach meant that the sample was limited in several ways. First, all participants were industry professionals, which is a subset of all software practitioners and might not be representative. Second, the participants were selected by managers, who might have their agenda in what employees to select. Third, the companies belonged to the subset of companies that were both sufficiently interested in this study and could allocate resources (i.e., subjects).

However, our results show that the effects of different professional characteristics, such as programming language and role, were limited. This could indicate that the study is less susceptible to convenience sampling than otherwise. Further, the 40 participants had a wide variety of professional backgrounds and were employed at twelve different companies and one government agency.

Along the same line, the generalizability of the results of the study is threatened by demographic factors. Due to various constraints (including financial), all partaking entities had offices in Sweden. While the study was conducted in several parts of the country, Sweden is culturally distinguished in terms of secular-rational and self-expression values (Inglehart and Welzel 2010). That said, there was diversity in, e.g., ethnicity among the participants, but such data was not collected to protect confidentiality. For the same reason, many aspects of the participants' demographic profiles were not investigated.

Finally, it is important to recognize that this study was an exploratory one, and not comprehensive. Hence, the quantitative findings should *not* be understood as applying to all design smells nor all instances of the selected design smells. What the data demonstrates is that—even in the context of small, isolated code examples—software practitioners' affective states can change in the presence of certain design smells.

7.4 Conclusion Validity

As far as we can tell, no previous studies have investigated how DTD relates to affects. Consequently, the findings of this study cannot be compared and contrasted with the findings of others. Instead, they must be evaluated in isolation and are, therefore, more susceptible to incorrect inferences and conclusions.

Three triangulation techniques (Miller 2008) were adopted to combat these threats. First, the data were triangulated in the sense that the sessions were spread out over four weeks, and the participants were employed at different entities. Second, researcher triangulation was achieved as two researchers took part in all data gathering and interpretation. Third, methodological triangulation was used, as data were collected through an experiment, a questionnaire, and interviews.

8 Conclusion

Fully understanding the impact of technical debt (TD) in the code base is a crucial challenge for researchers and practitioners alike. Although previous research has highlighted how TD can impact developers' morale, there is scarce evidence on how specific technical debt issues impact practitioners' affective states. Even more challenging is finding evidence related to design and architectural debt.

With our study, encompassing a quantitative data collection and analysis supported by additional qualitative insights from the participants, we offer a first detailed look into how the presence of design debt issues affect software practitioners' affective state.

The results show that five different smells have different impacts. Even when present in a small example, cyclic dependencies clearly and negatively affect software practitioners' affects. Simultaneously, missing encapsulation seems to be a more straightforward issue to deal with (although mildly affecting the practitioners' affects). Two issues related to hierarchy (missing hierarchy and rebellious hierarchy) seem to have a conflicting or no evident effect on the participants' affective state. In contrast, surprisingly, the presence of the broken modularization issue seems to have a positive impact on practitioners' affects.

These results imply that these different TDs need to be treated differently and that studying their impact on the practitioners' affective states helps to understand their overall impact (interest) and consequently how to prioritize them in practice. More studies with additional TD should be studied in a similar way as it was done in this study, so to provide a comprehensive catalog of the smells and their impact.

From our qualitative findings, it seems that practitioners undergo different levels of maturity in how they deal with TD. First, they might naively tend to avoid it (*Procrastination*), then they tend to build a quality-heavy mindset (mostly, however, by blaming others for the presence of TD, i.e., *Elitism*). Finally, they reach a higher level of maturity when a constructive mindset promotes high-quality code (*Compensation*). Also, practitioners seem to be affected negatively when they forecast TD, especially with *Apprehension* related to the future negative impact generated by TD, and by the inherent difficulty in identifying TD and predicting its consequences (TD as *Indeterminable* items).

Finally, we investigated whether participants' background covariates played a role, and we found partly how experience seems to act as a sort of amplifier for the participants' feelings, probably due to repeated encounters with TD and to the different maturity, acquired with more experience, in dealing with TD.

In summary, only some of the known issues highlighted in the literature seem to affect practitioners' feelings. At the same time, we find that dealing with TD is stressful and might require a fair amount of experience in the team to be handled constructively.

This topic remains mostly uncharted, and presents many opportunities for future work. A singular study is insufficient to build a solid theory, but we encourage others to replicate our experiment under similar or different settings, e.g., design smells, TD type, or cultures. Two particularly interesting investigations would be using industry code examples and situations that simulate time pressure.

Acknowledgements We want to thank all the participating companies, the individual participants in our study, and all the students who participated in our pilot studies. The computations were enabled by resources provided by the Swedish National Infrastructure for Computing (SNIC), partially funded by the Swedish Research Council through grant agreement no. 2018–05973.

Funding Open access funding provided by Chalmers University of Technology.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Al-Mutawa HA, Dietrich J, Marsland S, McCartin C (2014) On the shape of circular dependencies in Java programs. In: 23rd Australian software engineering conference ASWEC. IEEE Computer Society, pp 48–57. <https://doi.org/10.1109/ASWEC.2014.15>
- Alves NSR, Ribeiro LF, Caires V, Mendes TS, Spínola RO (2014) Towards an ontology of terms on technical debt. In: Sixth international workshop on managing technical debt. IEEE Computer Society, pp 1–7. <https://doi.org/10.1109/MTD.2014.9>
- Alves NSR, Mendes TS, de Mendonça Neto MG, Spínola RO, Shull F, Seaman CB (2016) Identification and management of technical debt: A systematic mapping study. *Inf Softw Technol* 70:100–121. <https://doi.org/10.1016/j.infsof.2015.10.008>
- Ampatzoglou A, Ampatzoglou A, Chatzigeorgiou A, Avgeriou P (2015) The financial aspect of managing technical debt: A systematic literature review. *Inf Softw Technol* 64:52–73. <https://doi.org/10.1016/j.infsof.2015.04.001>
- Avgeriou P, Kruchten P, Ozkaya I, Seaman C (2016) Managing technical debt in software engineering. *Dagstuhl Rep* 6(4):110–138. <https://doi.org/10.4230/DagRep.6.4.110>
- Besker T, Martini A, Bosch J (2017) The pricey bill of technical debt: When and by whom will it be paid. In: International conference on software maintenance and evolution ICSME. IEEE Computer Society, pp 13–23. <https://doi.org/10.1109/ICSME.2017.42>
- Besker T, Martini A, Bosch J (2018a) Managing architectural technical debt: A unified model and systematic literature review. *J Syst Softw* 135:1–16. <https://doi.org/10.1016/j.jss.2017.09.025>
- Besker T, Martini A, Lokuge RE, Blincoc K, Bosch J (2018b) Embracing technical debt, from a startup company perspective. In: 2018 International conference on software maintenance and evolution ICSME. IEEE Computer Society, pp 415–425. <https://doi.org/10.1109/ICSME.2018.00051>
- Besker T, Martini A, Bosch J (2019) Software developer productivity loss due to technical debt - A replication and extension study examining developers' development work. *J Syst Softw* 156:41–61. <https://doi.org/10.1016/j.jss.2019.06.004>
- Besker T, Ghanbari H, Martini A, Bosch J (2020) The influence of technical debt on software developer morale. *J Syst Softw* 167:110586. <https://doi.org/10.1016/j.jss.2020.110586>
- Betella A, Verschure PFMJ (2016) The affective slider: A digital self-assessment scale for the measurement of human emotions. *PLoS One* 11(2):e0148037. <https://doi.org/10.1371/journal.pone.0148037>
- Boehm BW, Papaccio PN (1988) Understanding and controlling software costs. *IEEE Trans Softw Eng* 14(10):1462–1477. <https://doi.org/10.1109/32.6191>
- Bradley MM, Lang PJ (1994) Measuring emotion: The self-assessment manikin and the semantic differential. *J Behav Ther Exp Psychiatry* 25(1):49–59. [https://doi.org/10.1016/0005-7916\(94\)90063-9](https://doi.org/10.1016/0005-7916(94)90063-9)
- Braun V, Clarke V (2006) Using thematic analysis in psychology. *Qual Res Psychol* 3(2):77–101. <https://doi.org/10.1191/1478088706qp063oa>
- Bürkner PC (2017) brms: An R package for Bayesian multilevel models using Stan. *J Stat Softw* 80:1–28. <https://doi.org/10.18637/jss.v080.i01>
- Bürkner PC (2018) Advanced Bayesian multilevel modeling with the R package brms. *R J* 10(1):395–411. <https://doi.org/10.32614/RJ-2018-017>
- Bürkner PC, Vuorre M (2019) Ordinal regression models in psychology: A tutorial. *Adv Methods Pract Psychol Sci* 2(1):77–101. <https://doi.org/10.1177/2515245918823199>
- Colomo Palacios R, Hernández-López A, García-Crespo Á, Soto-Acosta P (2010) A study of emotions in requirements engineering. In: Lytras MD, de Pablos PO, Ziderman A, Roulstone A, Maurer HA, Imber JB (eds) Organizational, business, and technological aspects of the knowledge society—third world summit on the knowledge society WSKS, Springer, Communications in Computer and Information Science, vol 112. pp 1–7. https://doi.org/10.1007/978-3-642-16324-1_1
- Cruz SSJO, da Silva FQB, Capretz LF (2015) Forty years of research on personality in software engineering: A mapping study. *Comput Hum Behav* 46:94–113. <https://doi.org/10.1016/j.chb.2014.12.008>
- Cruzes DS, Dybå T (2011) Recommended steps for thematic synthesis in software engineering. In: Proceedings of the 5th International symposium on empirical software engineering and measurement ESEM. IEEE Computer Society, pp 275–284. <https://doi.org/10.1109/ESEM.2011.36>
- Cunningham W (1992) The WyCash portfolio management system. *ACM SIGPLAN OOPS Messenger* 4(2):29–30. <https://doi.org/10.1145/157710.157715>
- Ernst NA, Bellomo S, Ozkaya I, Nord RL, Gorton I (2015) Measure it? Manage it? Ignore it? Software practitioners and technical debt. In: Nitto ED, Harman M, Heymans P (eds) Proceedings of the 10th joint meeting on foundations of software engineering ESEC/ FSE. ACM, 50–60. <https://doi.org/10.1145/2786805.2786848>

- Fagerholm F, Ikonen M, Kettunen P, Münch J, Roto V, Abrahamsson P (2015) Performance alignment work: How software developers experience the continuous adaptation of team performance in Lean and Agile environments. *Inf Softw Technol* 64:132–147. <https://doi.org/10.1016/j.infsof.2015.01.010>
- Feldt R, Angelis L, Torkar R, Samuelsson M (2010) Links between the personalities, views and attitudes of software engineers. *Inf Softw Technol* 52(6):611–624. <https://doi.org/10.1016/j.infsof.2010.01.001>
- Fernández-Sánchez C, Garbajosa J, Yagüe A, Pérez J (2017) Identification and analysis of the elements required to manage technical debt by means of a systematic mapping study. *J Syst Softw* 124:22–38. <https://doi.org/10.1016/j.jss.2016.10.018>
- Fontana FA, Pigazzini I, Roveda R, Tamburri DA, Zanoni M, Nitto ED (2017) Arcan: A tool for architectural smells detection. In: International conference on software architecture workshops ICSA. IEEE Computer Society, pp 282–285. <https://doi.org/10.1109/ICSAW.2017.16>
- Furia CA, Feldt R, Torkar R (2019) Bayesian data analysis in empirical software engineering research. *IEEE Trans Softw Eng* :1–1. accepted for publication
- Ganesh SG, Sharma T, Suryanarayana G (2013) Towards a principle-based classification of structural design smells. *J Object Technol* 12(2):1–29. <https://doi.org/10.5381/jot.2013.12.2.a1>
- Garcia J, Popescu D, Edwards G, Medvidovic N (2009) Toward a catalogue of architectural bad smells. In: Mirandola R, Gorton I, Hofmeister C (eds) 5th International conference on the quality of software architectures QoSA, Springer, Lecture Notes in Computer Science, vol 5581. pp 146–162. https://doi.org/10.1007/978-3-642-02351-4_10
- Gelman A (2018) The failure of null hypothesis significance testing when studying incremental changes, and what to do about it. *Pers Soc Psychol Bull* 44(1):16–23. <https://doi.org/10.1177/0146167217729162>
- Gelman A, Tuerlinckx F (2000) Type S error rate for classical and Bayesian single and multiple comparison procedures. *Comput Stat* 15:373–390. <https://doi.org/10.1007/s001800000040>
- Gelman A, Hill J, Yajima M (2012) Why we (usually) don't have to worry about multiple comparisons. *J Res Educ Effect* 5(2):189–211. <https://doi.org/10.1080/19345747.2011.618213>
- Goodman SN (1999a) Toward evidence-based medical statistics. 1: The p value fallacy. *Ann Intern Med* 130(12):995–1004. <https://doi.org/10.7326/0003-4819-130-12-199906150-00008>
- Goodman SN (1999b) Toward evidence-based medical statistics. 2: The Bayes factor. *Ann Intern Med* 130(12):1005–1013. <https://doi.org/10.7326/0003-4819-130-12-199906150-00019>
- Graziotin D, Wang X, Abrahamsson P (2014) Happy software developers solve problems better: Psychological measurements in empirical software engineering. *PeerJ* 2:e289. <https://doi.org/10.7717/peerj.289>
- Graziotin D, Wang X, Abrahamsson P (2015a) The affect of software developers: common misconceptions and measurements. In: 2015 IEEE/ACM 8th International workshop on cooperative and human aspects of software engineering. IEEE, pp 123–124
- Graziotin D, Wang X, Abrahamsson P (2015b) Do feelings matter? On the correlation of affects and the self-assessed productivity in software engineering. *J Softw Evol Process* 27(7):467–487. <https://doi.org/10.1002/smr.1673>
- Graziotin D, Wang X, Abrahamsson P (2015c) Understanding the affect of developers: Theoretical background and guidelines for psychoempirical software engineering. In: Hammouda I, Sillitti A (eds) 7th International workshop on social software engineering SSE. ACM, pp 25–32. <https://doi.org/10.1145/2804381.2804386>
- Graziotin D, Fagerholm F, Wang X, Abrahamsson P (2017) On the unhappiness of software developers. In: Mendes E, Counsell S, Petersen K (eds) 21st International conference on evaluation and assessment in software engineering EASE. ACM, pp 324–333. <https://doi.org/10.1145/3084226.3084242>
- Graziotin D, Fagerholm F, Wang X, Abrahamsson P (2018) What happens when software developers are (un)happy. *J Syst Softw* 140:32–47. <https://doi.org/10.1016/j.jss.2018.02.041>
- Gren L, Torkar R, Feldt R (2017) Group development and group maturity when building agile teams: A qualitative and quantitative investigation at eight large companies. *J Syst Softw* 124:104–119. <https://doi.org/10.1016/j.jss.2016.11.024>
- Guo Y, Seaman CB, Gomes R, Cavalcanti ALO, Tonin G, da Silva FQB, de Medeiros Santos AL, de Siebra C (2011) Tracking technical debt—An exploratory case study. In: 27th International conference on software maintenance ICSM. IEEE Computer Society, pp 528–531. <https://doi.org/10.1109/ICSM.2011.6080824>
- Inglehart R, Welzel C (2010) Changing mass priorities: The link between modernization and democracy. *Perspect Polit* 8(2):551–567. <https://doi.org/10.1017/S1537592710001258>
- Khan IA, Brinkman WP, Hierons RM (2011) Do moods affect programmers' debug performance? *Cognit Technol Work* 13(4):245–258. <https://doi.org/10.1007/s10111-010-0164-1>

- Kruschke JK (2010) What to believe: Bayesian methods for data analysis. *Trends Cognit Sci* 14(7):293–300. <https://doi.org/10.1016/j.tics.2010.05.001>
- Lang PJ (1980) Behavioral treatment and bio-behavioral assessment: Computer applications. *Technology in mental health care delivery systems*. pp 119–137
- Lang PJ, Bradley MM, Cuthbert BN (1997) International affective picture system (IAPS): Technical manual and affective ratings. *NIMH Center for the Study of Emotion and Attention* 1:39–58
- Leek J, McShane BB, Gelman A, Colquhoun D, Nuijten MB, Goodman SN (2017) Five ways to fix statistics. *Nature* :557–559. <https://doi.org/10.1038/d41586-017-07522-z>
- Lenberg P, Feldt R, Wallgren LG (2015) Behavioral software engineering: A definition and systematic literature review. *J Syst Softw* 107:15–37. <https://doi.org/10.1016/j.jss.2015.04.084>
- Li Z, Avgeriou P, Liang P (2015) A systematic mapping study on technical debt and its management. *J Syst Softw* 101:193–220. <https://doi.org/10.1016/j.jss.2014.12.027>
- Lim E, Taksande N, Seaman C (2012) A balancing act: What software practitioners have to say about technical debt. *IEEE Softw* 29(6):22–27. <https://doi.org/10.1109/MS.2012.130>
- Martini A, Bosch J (2017) On the interest of architectural technical debt: Uncovering the contagious debt phenomenon. *J Softw Evol Process* 29(10):e1877. <https://doi.org/10.1002/smr.1877>
- Martini A, Besker T, Bosch J (2018a) Technical debt tracking: Current state of practice: A survey and multiple case study in 15 large organizations. *Sci Comput Program* 163:42–61. <https://doi.org/10.1016/j.scico.2018.03.007>
- Martini A, Fontana FA, Biaggi A, Roveda R (2018b) Identifying and prioritizing architectural debt through architectural smells: A case study in a large software company. In: Cuesta CE, Garlan D, Pérez J (eds) 12th European conference on software architecture ECSA, Springer, Lecture Notes in Computer Science, vol 11048. pp 320–335. https://doi.org/10.1007/978-3-030-00761-4_21
- Miller J (2008) Triangulation as a basis for knowledge discovery in software engineering. *Empir Softw Eng* 13(2):223–228. <https://doi.org/10.1007/s10664-008-9063-y>
- Morris JD (1995) Observations: SAM: The self-assessment manikin: An efficient cross-cultural measurement of emotional response. *J Advert Res* 35:63–68
- Morris JD, Woo C, Geason JA, Kim J (2002) The power of affect: Predicting intention. *J Advert Res* 42(3):7–17. <https://doi.org/10.2501/JAR-42-3-7-17>
- Peterson C, Park N, Sweeney PJ (2008) Group well-being: morale from a positive psychology perspective. *Appl Psychol* 57:19–36
- R Core Team (2020) R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>
- Runeson P, Höst M (2009) Guidelines for conducting and reporting case study research in software engineering. *Empir Softw Eng* 14(2):131–164. <https://doi.org/10.1007/s10664-008-9102-8>
- Russell JA, Mehrabian A (1977) Evidence for a three-factor theory of emotions. *J Res Pers* 11(3):273–294. [https://doi.org/10.1016/0092-6566\(77\)90037-X](https://doi.org/10.1016/0092-6566(77)90037-X)
- Sharma T, Spinellis D (2018) A survey on software smells. *J Syst Softw* 138:158–173. <https://doi.org/10.1016/j.jss.2017.12.034>
- Spínola RO, Vetrò A, Zazworka N, Seaman C, Shull F (2013) Investigating technical debt folklore: Shedding some light on technical debt opinion. In: 4th International workshop on managing technical debt MTD. pp 1–7. <https://doi.org/10.1109/MTD.2013.6608671>
- Stan Development Team (2020) RStan: The R interface to Stan. <http://mc-stan.org/>, r package version 2.21.1
- Suryanarayana G, Samarthyam G, Sharma T (2014) Refactoring for software design smells: Managing technical debt. Burlington, Morgan Kaufmann
- Tamburri DA, Kruchten P, Lago P, van Vliet H (2013) What is social debt in software engineering. In: 6th International workshop on cooperative and human aspects of software engineering CHASE. IEEE Computer Society, pp 93–96. <https://doi.org/10.1109/CHASE.2013.6614739>
- Tom E, Aurum A, Vidgen R (2012) An exploration of technical debt. *J Syst Softw* 86(6):1498–1516. <https://doi.org/10.1016/j.jss.2012.12.052>
- Vehtari A, Gelman A, Gabry J (2017) Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Stat Comput* 27:1413–1432. <https://doi.org/10.1007/s11222-016-9696-4>
- Weller SC, Vickers B, Bernard HR, Blackburn AM, Borgatti S, Gravlee CC, Johnson JC (2018) Open-ended interview questions and saturation. *PLoS ONE* 13(6):1–18. <https://doi.org/10.1371/journal.pone.0198606>
- Wohlin C, Runeson P, Höst M, Ohlsson MC, Regnell B, Wesslén A (2012) Experimentation in software engineering. Springer Science & Business Media, Berlin

Yli-Huumo J, Maglyas A, Smolander K (2014) The sources and approaches to management of technical debt: A case study of two product lines in a middle-size finnish software company. In: Jedlitschka A, Kuvaja P, Kuhrmann M, Männistö T, Münch J, Raatikainen M (eds) 15th International conference on product-focused software process improvement PROFES, Springer, Lecture Notes in Computer Science, vol 8892. pp 93–107. https://doi.org/10.1007/978-3-319-13835-0_7

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.